

BDS Exam Preparations.....	2
Parallel Programming.....	2
Amdahl's Law.....	2
Gustafson's Law.....	2
Realistic Data Pipeline.....	5
Estimating Power Efficiency of Dual Core CPU.....	5
OpenMP.....	5
Data Sharing:.....	6
MPI Programs.....	7
Parallel I/O Systems.....	9
Fault Tolerance.....	9
Disk Level.....	9
RAID.....	9
Node Level.....	10
Google File System (GFS).....	11
MapReduce.....	12
In Memory Computing.....	14
RDD (Resilient Distributed Datasets).....	15
Spark Programming.....	15
Graph Analysis.....	17
Asynchronous processing.....	17
Graph partitioning on power law graphs and GAS model.....	17
Out of core graph processing.....	18
GRAPHCHI – Parallel Sliding Window.....	19
XSTREAM – Edge Centric.....	19
Streaming Partitions.....	19
GRIDGRAPH – 2D Partition.....	20
Top Down vs. Bottom Up.....	20
GEMINI – Distributed push/pull.....	21
Shentu - Designed for Extreme Scale Graph.....	21
Graph Mining.....	22
Stream Processing.....	22
Actor Model.....	22
Stream Processing – STORM.....	23
Core Concepts and Architecture.....	23
SparkStream.....	24
Comparison of different Streaming Models.....	24
Machine Learning.....	24

Halide.....	25
Graph Optimization.....	25
Distributed Training.....	25
Memory Consumption.....	26
AdaPipe.....	26
Blockchain.....	27
Conflux.....	28
Computer Systems.....	29
Why is it difficult to build a system?.....	29
Common Problems of Systems.....	29
Coping with Complexity.....	29
Prep. Exam.....	30

BDS Exam Preparations

Challenges of Big Data Processing:

- Volume → More memory and I/O bandwidth
- Velocity → Fast Processing
- Variety → More computing paradigms

Parallel Programming

Challenges:

- Parallelism identification
- Load balancing
- Proper synchronization

Amdahl's Law

If P is the proportion of a program that can be made parallel, the limit of speedup of whole program is bounded at **Speed-up = 1/(1-P)**

Gustafson's Law

Amdahl's law assumes fixed problem size

If problem size can increase as the compute power increases, the speedup that can be achieved is much larger.

S = a(n) + p(1-a(n)))

n = problem size

p = number of processes

a(n) = serial portion at problem size n

So S → p when a(n) = 0

⇒ Larger machines are still useful

⇒ Many small chips can deliver good performance when the problem size is big

Problem: Some applications don't have bigger data set

Real world is mix of Amdahl's law and Gustafson's law

Posix Threads (Pthreads) Interface

-

Pthreads:

Standard interface for ~60 functions that

manipulate threads from C programs

–Creating and reaping threads

- `pthread_create()`

- `pthread_join()`

–Determining your thread ID

- `pthread_self()`

–Terminating threads

- `pthread_cancel()`

- `pthread_exit()`

•`exit()` [terminates all threads] ,`RET` [terminates current thread]

–Synchronizing access to shared variables

- `pthread_mutex_init`

- `pthread_mutex_[un]lock`

- `pthread_cond_init`

- `pthread_cond_[timed]wait`

Semaphores:

Non-negative global integer synchronization variables

Manipulated by P and V operations:

- P(s): [`while (s == 0) wait(); s--;`]

- Dutch for "Proberen" (test)

- V(s): [`s++;`]

- Dutch for "Verhogen" (increment)

OS kernel guarantees that operations between brackets [] are executed indivisibly

- Only one P or V operation at a time can modify s.
- When while loop in P terminates, only that P can decrement s

```
void *sum_thread(void *vargp)
{
    long sum;
    int myid = (long) vargp;
    int start = myid * sample_num / 4;
    int end = (myid+1) * sample_num / 4 - 1;

    sum = 0;
    for (int i=start; i<=end; i++)
        sum += x[i];

    sem_wait(&num_mutex);
    global_sum += sum;
    sem_post(&num_mutex);

    return NULL;
}
```

Issues with Pthreads:

Complex for data processing

- Pthread is not always portable(e.g. Windows)
- Code to calculate addresses and loop bounds
- May not be easy to specify scheduling schemes

```
void *sum_thread(void *);

volatile long global_sum = 0; /* global */

main()
{
    ...
    for (int i = 0; i < N; i++)
        pthread_create(&tid[i], NULL, sum_thread, (void*)i);

    for (int i = 0; i < N; i++)
        pthread_join(tid[i], NULL);
}
```

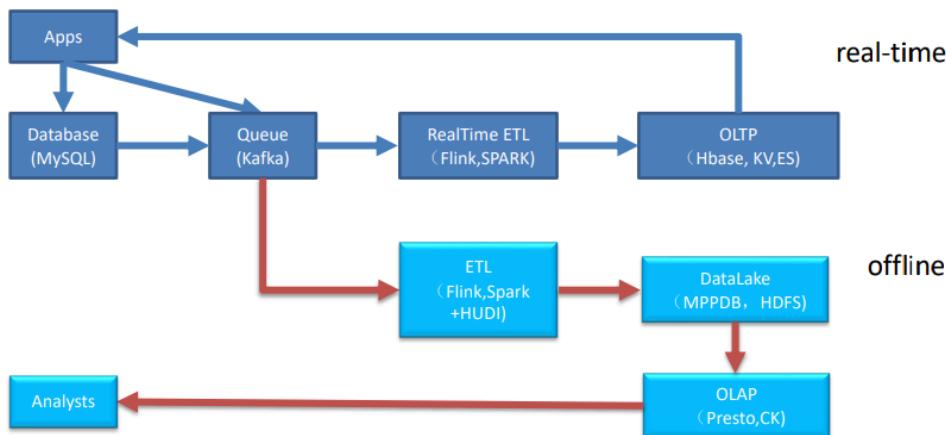
```
void *sum_thread(void *);

volatile long global_sum = 0; /* global */

main()
{
    ...
    for (int i = 0; i < N; i++)
        pthread_create(&tid[i], NULL, sum_thread, (void*)i);

    for (int i = 0; i < N; i++)
        pthread_join(tid[i], NULL);
}
```

Realistic Data Pipeline



Estimating Power Efficiency of Dual Core CPU

Single Core: 1 solution/s

$$P = cf_1^2$$

$$P/2 = cf_2^2$$

$$F_2 = f_1/(2)^{1/2}$$

$$\text{Dual Core Throughput} = 2*f_2 = 2^{1/2} * f_1 = 2^{1/2} \text{ solutions/s}$$

⇒ With perfect parallelism

OpenMP

“Standard” API for defining multi-threaded sharedmemory programs

– Allow a programmer to **separate a program into serial regions and parallel regions**, rather than T concurrently-executing threads.

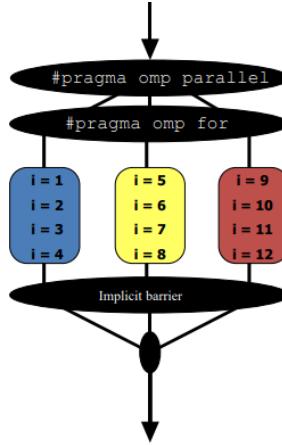
– **Hide stack management**

– Provide synchronization constructs

⇒ **ONLY SINGLE MACHINES**

- Threads are assigned an independent set of iterations
- Threads must wait at the end of work-sharing construct

```
#pragma omp parallel for
for(i = 1, i < 13, i++)
    c[i] = a[i] + b[i]
```



Data Sharing:

- Parallel programs often employ two types of data
 - Shared data, visible to all threads, similarly named
 - Private data, visible to a single thread (often stack-allocated)
- pthreads:
 - Global-scoped variables are shared
 - Stack-allocated variables are private
- OpenMP:
 - **shared** variables are shared
 - **private** variables are private
 - Rules for default

```
int bigdata[1024];

void* foo(void* bar) {
    int tid;

    #pragma omp parallel \
        shared ( bigdata ) \
        private ( tid )
    {
        /* Calc. here */
    }
}

int i = 0;
int n = 10;
int a = 7;

#pragma omp parallel for
for (i = 0; i < n; i++)
{
    int b = a + i;
    ...
}
```

n and a are shared. i and b are private.

Firstprivate

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

    #pragma omp parallel for firstprivate(sum)
    for (int i = 1; i <= 10; i++)
    {
        sum += i;
        printf("In iteration %d, sum is %d\n", i, sum);
    }

    printf("After the loop, sum is %d\n", sum);
}
```

- Initialize private variables from the master threads

lastprivate

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

    #pragma omp parallel for firstprivate(sum) lastprivate(sum)
    for (int i = 1; i <= 10; i++)
    {
        sum += i;
        printf("In iteration %d, sum is %d\n", i, sum);
    }

    printf("After the loop, sum is %d\n", sum);
}
```

- Pass the value of a private from the **last iteration** to global variable

OpenMP Critical Construct

- `#pragma omp critical [(lock_name)]`
- Defines a critical region on a structured block

```
Threads wait their turn -at
a time, only one calls
consum() thereby
protecting RES from race
conditions
```

```
float RES;
#pragma omp parallel
{ float B;
#pragma omp for private(B)
for(int i=0; i<niters; i++){
    B = big_job(i);
#pragma omp critical
    {
        consum (B, RES);
    }
}
```

Sum of an Array

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

#pragma omp parallel for reduction(+:sum)
for (int i = 1; i <= 10; i++)
{
    sum += i;
    printf("In iteration %d, sum is %d\n", i, sum);
}

printf("After the loop, sum is %d\n", sum);
```

- The `schedule clause` affects how loop iterations are mapped onto threads

```
schedule (static [,chunk])
    • Blocks of iterations of size "chunk" to threads
    • Round robin distribution
    • Default=N/t

schedule (dynamic [,chunk])
    • Threads grab "chunk" iterations
    • When done with iterations, thread requests next set
    • Default=1

schedule (guided [,chunk])
    • Dynamic schedule starting with large block
    • Size of the blocks shrink; no smaller than "chunk"
    • Default=1

schedule (runtime)
    • OMP_SCHEDULE
```

Schedule Clause	When To Use
STATIC	Predictable and similar work per iteration
DYNAMIC	Unpredictable, highly variable work per iteration
GUIDED	Special case of dynamic to reduce scheduling overhead

MPI Programs

```
#include "mpi.h"
#include <stdio.h>

int main( int argc, char *argv[] )
{
    int rank, size;
    MPI_Init( &argc, &argv );
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
    MPI_Comm_size( MPI_COMM_WORLD, &size );
    printf( "I am %d of %d\n", rank, size );
    MPI_Finalize();
    return 0;
}
```

```
Process P:
send(request1,32,Q)

Process R:
send(request2, 32, Q)

Process Q:
while (true) {
    recv(received_request, 32, Any_Process);
    process received_request;
}
```

```
Process P:
send(request1, 32, Q, tag1)

Process R:
send(request2, 32, Q, tag2)

Process Q:
while (true){
    recv(received_request, 32, Any_Process, Any_Tag, Status);
    if (Status.MPI_TAG==tag1) process received_request in one way;
    if (Status.MPI_TAG==tag2) process received_request in another way;
}
```

Tags can also serve as "message type", although we should not confuse it with MPI data types

```
MPI_Init(&argc,&argv);
MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
MPI_Comm_rank(MPI_COMM_WORLD,&myid);

tag=1234;source=0;destination=1;count=1;

if(myid == source){
    buffer=5678;
    MPI_Send(&buffer,count,MPI_INT,destination,tag,MPI_COMM_WORLD);
    printf("processor %d sent %d\n",myid,buffer);
}
if(myid == destination){
    MPI_Recv(&buffer,count,MPI_INT,source,tag,MPI_COMM_WORLD,&status);
    printf("processor %d got %d\n",myid,buffer);
}

MPI_Finalize();
```

MPI Basic (Blocking) Send/Receive

`MPI_SEND (start, count, datatype, dest, tag, comm)`
`MPI_RECV(start, count, datatype, source, tag, comm, status)`

Why Datatypes? ⇒ Support communication between processes on machines with different memory representations

Unsafe message passing ⇒ Depending on the availability of system buffers:

- Order the operations more carefully:

Process 0	Process 1
Send(1)	Recv(0)
Recv(1)	Send(0)

- Use non-blocking operations:

Process 0	Process 1	Process 0	Process 1
Send(1)	Send(0)	Isend(1)	Irecv(0)
Recv(1)	Recv(0)	Irecv(1)	Send(0)

Better:

```
if (rank == 0) {
    MPI_Isend(..., 1, tag, MPI_COMM_WORLD, &req);
    MPI_Recv(..., 1, tag, MPI_COMM_WORLD, &status);
    MPI_Wait(&req, &status);
} else if (rank == 1) {
    MPI_Recv(..., 0, tag, MPI_COMM_WORLD, &status);
    MPI_Send(..., 0, tag, MPI_COMM_WORLD);
}
```

Just six functions:

– MPI_INIT – MPI_FINALIZE – MPI_COMM_SIZE – MPI_COMM_RANK – MPI_SEND – MPI_RECV

Collective Operations in MPI:

- Collective operations are called by all processes in a communicator.
- **MPI_BCAST** distributes data from one process (the root) to all others in a communicator.
- **MPI_REDUCE** combines data from all processes in communicator and returns it to one process.
- In many numerical algorithms, **SEND/RECEIVE** can be replaced by **BCAST/REDUCE**, improving both simplicity and efficiency.

Problem with MPI: Availability

What if a process crashes?

→ If app runs 20 days or longer

 → Restart doesn't help

Checkpointing is required!

What are the obstacles of supporting heterogeneous node and fault tolerance

- Knowing N at the beginning
- Static partition of the workload

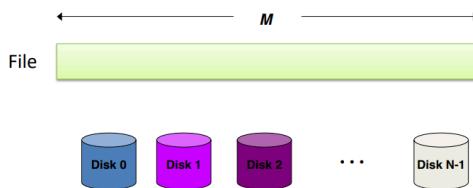
```
#include "mpi.h"
#include <math.h>
int main(int argc, char *argv[])
{
    int done = 0, n, myid, numprocs, i, rc;
    double PI25DT = 3.141592653589793238462643;
    double mypi, pi, h, sum, x, a;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Cart_create(MPI_COMM_WORLD, 1, numprocs, 1, 1, &rc);
    h = 1.0 / (double) n;
    sum = 0.0;
    for (i = myid + 1; i <= n; i += numprocs) {
        x = h * ((double)i - 0.5);
        sum += 4.0 / (1.0 + x*x);
    }
    MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,
               MPI_COMM_WORLD);
    if (myid == 0)
        printf("pi is approximately %.16f, Error is %.16f\n",
               pi, fabs(pi - PI25DT));
}
MPI_Finalize();
return 0;
}
```

Parallel I/O Systems

Fault Tolerance

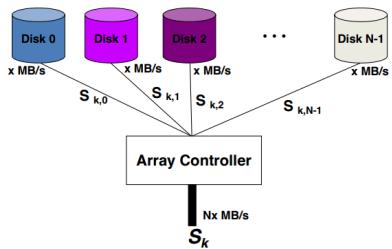
Disk Level

- Given a large file of M bytes file, and N disks
 - Find a way to distribute the file into disks.
- Objective: Read a long sequential part of the file as fast as possible



Performance Benefit

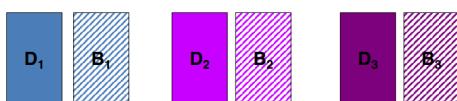
- Sequential read or write of large file*
 - application (or I/O buffer cache) reads in multiples of S bytes
 - controller performs parallel access of N disks
 - aggregate bandwidth is N times individual disk bandwidth
 - (assumes that disk is the bottleneck)



be useful

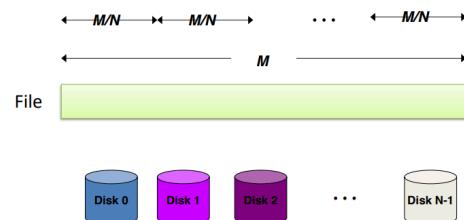
→ “Redundant Arrays of Independent Disks”

RAID Level 1



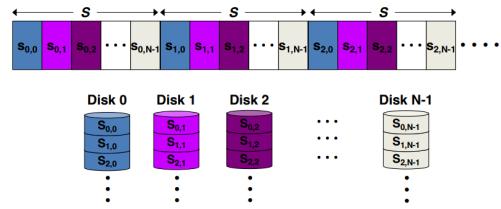
- Also known as “*mirroring*”
- To read a block:
 - read from either data disk or backup
- To write a block:
 - write both data and backup disks
- How many failure can tolerate? What is the cost for redundancy?

Solution 1: Chunk



Solution 2: Stripe

- Stripe* the data across an array of disks
 - many alternative striping strategies possible
- Example: consider a big file striped across N disks
 - stripe width* is S bytes
 - hence each *stripe unit* is S/N bytes
 - sequential read of S bytes at a time



RAID

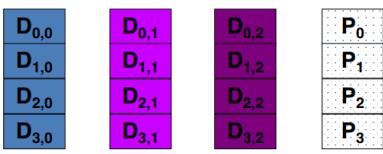
Large arrays without redundancy are too unreliable to

RAID Levels 2 & 3



- These are *bit-interleaved* schemes
- In Raid Level 2, P contains memory-style ECC
- In Raid Level 3, P contains simple parity
- Rarely used today
- Good for large data read/write, poor for small data transfer

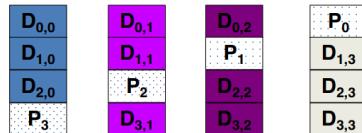
RAID Level 4



$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Block-interleaved parity*
- Wasted storage is small: one parity block for N data blocks
- Key problem:
 - parity disk becomes a hot spot
 - write access to parity disk on every write to any block

RAID Level 5



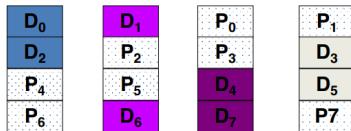
$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Rotated parity*
- Wastage is small: same as in Raid 4
- Parity update traffic is distributed across disks

	RAID-0	RAID-1	RAID-4	RAID-5
Capacity	N	$N/2$	$N - 1$	$N - 1$
Reliability	0	1 (for sure)	1	1
		$\frac{N}{2}$ (if lucky)		
Throughput				
Sequential Read	$N \cdot S$	$(N/2) \cdot S$	$(N - 1) \cdot S$	$(N - 1) \cdot S$
Sequential Write	$N \cdot S$	$(N/2) \cdot S$	$(N - 1) \cdot S$	$(N - 1) \cdot S$
Random Read	$N \cdot R$	$N \cdot R$	$(N - 1) \cdot R$	$N \cdot R$
Random Write	$N \cdot R$	$(N/2) \cdot R$	$\frac{1}{2} \cdot R$	$\frac{N}{4} \cdot R$
Latency				
Read	D	D	D	D
Write	D	D	$2D$	$2D$

Problem with RAID lvl 5 if failure rate not $>>$ recovery time.

RAID Level 6



RAID provide good reliability for storage in a single node

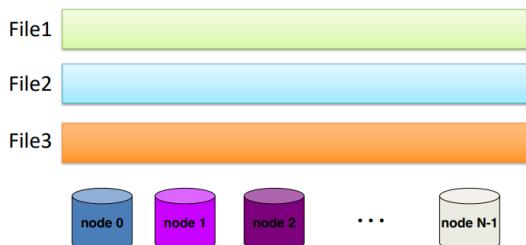
- But it **can not tolerate node/network failure**

- Tolerate 2 disk failures
- Read performance is good, while write performance is a little slower
- Parity update traffic is distributed across disks

Node Level

Problem Statement

- Given a set of large files and N machine nodes
 - Find a way to distribute the file into nodes.
- Objective:
 - Tolerate node failure
 - Read a long sequential part of the file as fast as possible
 - Write are mostly large append mode write

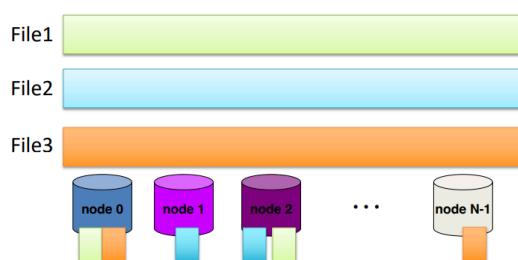


Duplicate Whole Files to Nodes

A node include all content of file.

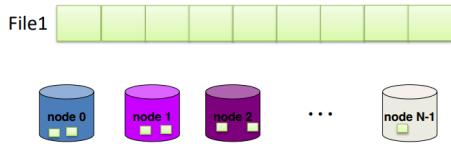
A file is duplicated in multiple nodes. (In this example:2)

Pros and cons of this strategy?



Duplicate chunks of files

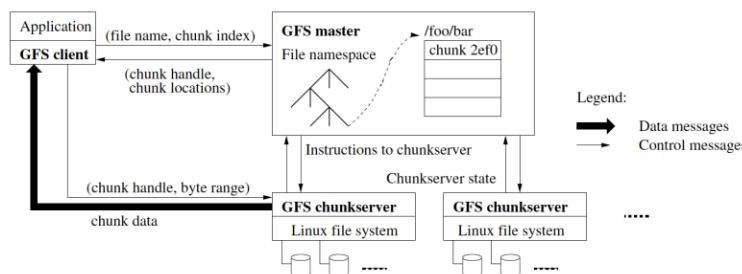
- A file is divided into fixed size chunks
- Each chunk is distributed and duplicated to a few nodes (usually 2-3)
- Pros and cons of this strategy?
 - Number of node failure tolerated?
 - Wasted space?
 - Space required to know the position of the chunk?



How to distribute 3 replicas on node/rack hierarchy?

- One replica on local node
- Second replica on a remote rack
- Third replica on same remote rack

Google File System (GFS)



Key GFS design decisions

- 3 replicas of data for fault tolerance
- Single master design for simplicity and performance
- Master does not in data paths
- Large chunk size for smaller meta data

Single master design

- Performance issue
 - Separate data path with control path, no data need to pass chunk server
 - Let's assume the system contains modest number of HUGE tiles, and the chunk size is large(e.g. 64MB)
 - Meta data for the system is not large
 - They can be put in memory of a single server, this greatly reduce the complexity of the system and could get good meta data performance.
- Single master design: Simple! With Availability and performance issues

Metadata in Memory

- The entire metadata is in main memory
- No demand paging of metadata

Data Correctness

Chunk server

A Block Server

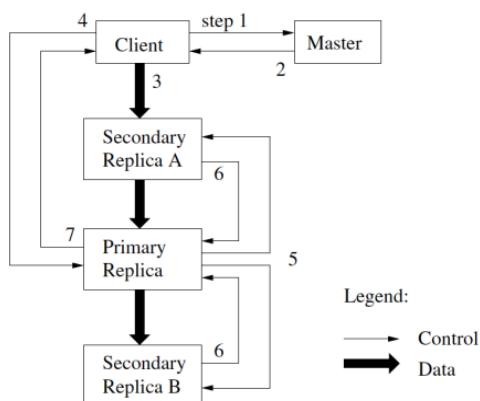
- Stores data in the local file system (e.g. ext3)
 - Each block of GFS data in a separate file in its local file system.
- Stores metadata of a block (e.g. CRC)
- Serves data and metadata to Clients

- Use Checksums to validate data
 - Use CRC32

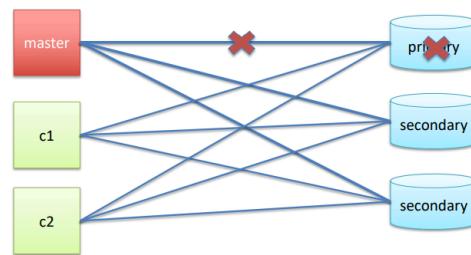
- File Creation
 - Client computes checksum per 512 bytes
 - Chunk server stores the checksum
- File access
 - Client retrieves the data and checksum from chunk server
 - If Validation fails, Client tries other replicas

GFS Write

A primary is need to control the order of the concurrent write

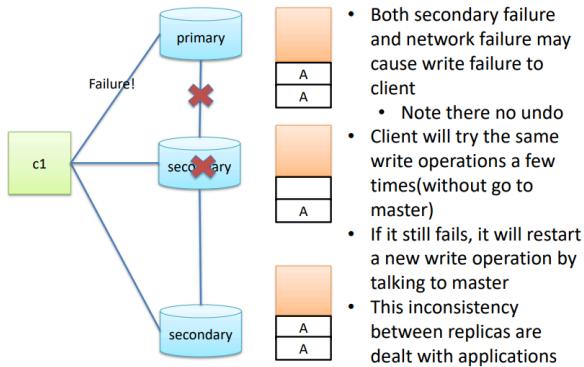


Why we need a lease?



A lease for primary is to deal with primary failure or network failure

What happens if secondary fails



Heartbeats

- Master sends heartbeat to the chunk servers
 - Once every 3 seconds
 - No heartbeat does not necessarily indicate chunk server failure
 - May be network failure

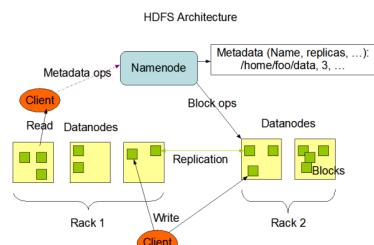
Rebalancer

- Goal: % disk full on chunk servers should be similar
 - Usually run when new chunk servers are added
 - Rebalancer is throttled to avoid network congestion

Replication Engine

- Master detects chunk server failures
 - Chooses new chunk server for new replicas
 - Balances disk usage
 - Balances communication traffic to chunk servers

HDFS Architecture



User Interface

- Commands for HDFS User:
 - hadoop dfs -mkdir /foodir
 - hadoop dfs -cat /foodir/myfile.txt
 - hadoop dfs -rm /foodir/myfile.txt
- Commands for HDFS Administrator
 - hadoop dfsadmin -report
 - hadoop dfsadmin -decommission datanodename
- Web Interface
 - http://host:port/dfshealth.jsp

Example: Count word occurrences (pseudo code)

```

map(String input_key, String input_value):
    // input_key: document name
    // input_value: document contents
    for each word w in input_value:
        EmitIntermediate(w, "1");

reduce(String output_key, Iterator intermediate_values):
    // output_key: a word
    // output_values: a list of counts
    int result = 0;
    for each v in intermediate_values:
        result += ParseInt(v);
    Emit(AsString(result));
  
```

MapReduce

- Doing a small piece of work, complete it and get the next piece
 - Input from the global file system, similar to the MPI_Bcast
 - Write a `map()` function that transfers the input piece into partial output piece
- No N is required to know at the beginning
 - Data specified by files in distributed file system
- Need a reduce support similar to MPI_Reduce
 - Write a user defined `reduce()` file

GFS design make an tradeoff

- Trade ideal consistency for system simplicity
- Clients may observe
 - Reading stale data
 - Duplicate append records
- For applications such as search engine
 - This may be tolerable
- Understanding GFS design as a tradeoff among
 - Consistency, fault-tolerance, performance, simplicity of design

Map output

Example

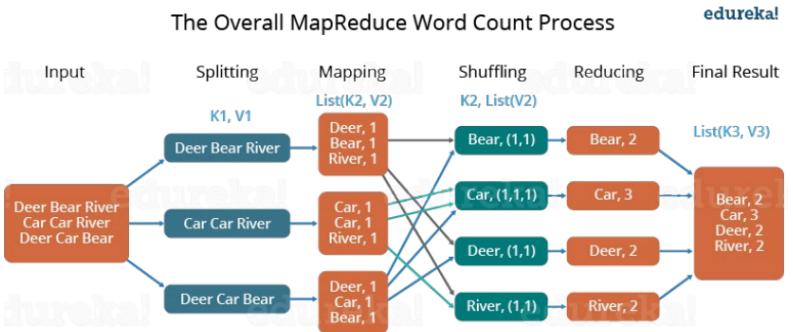
- Page 1: the weather is good
- Page 2: today is good
- Page 3: good weather is good.

Reduce Output

- Worker 1:
– (the 1)
- Worker 2:
– (is 3)
- Worker 3:
– (weather 2)
- Worker 4:
– (today 1)
- Worker 5:
– (good 4)

Parallelism

- map() functions run in parallel, creating different intermediate values from different input data sets
- reduce() functions also run in parallel, each working on a different output key
- All values are processed *independently*
- Bottleneck: reduce phase can't start until map phase is completely finished.



Why MapReduce is successful?

- Programmers only need to write serial code
 - Shared/private code in OpenMP
 - message orders in MPI
- Automatically parallel and distributed computing
- Automatic load-balancing and fault-tolerance

Limitations of MapReduce

- Only Map and Reduce
- Iterative MapReduce
 - Intermediate result on disk
 - Many I/O operations, poor performance
- Unable to support interactive query efficiently
 - No mechanism to share data in memory between MapReduce jobs

Reduce Input

- Worker 1:
– (the 1)
- Worker 2:
– (is 1), (is 1), (is 1)
- Worker 3:
– (weather 1), (weather 1)
- Worker 4:
– (today 1)
- Worker 5:
– (good 1), (good 1), (good 1), (good 1)

Transferring data with file system



- Iterative MapReduce is required when you need to "sum up" results more than 1 times
- The MapReduce model is stateless, causing performance problems for iterative MapReduce programs

Expressing an Iteration in MapReduce

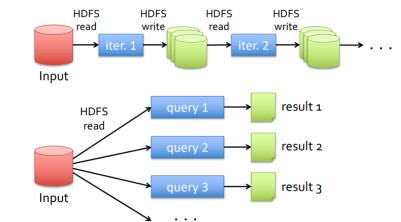
- for each point in the point set
 - find the closest mean i , assign the point to cluster i , forming new K clusters
- calculate new means for each cluster

```

map(point p):
  for each mean in K means
    find the closest mean M to this point
    emit_intermediate(index of M, p)

reduce(index of cluster i, list of points lp):
  calculate the new mean M' for points in lp
  emit(M')
  
```

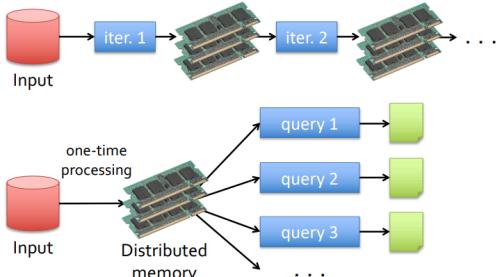
MapReduce Use File to Transfer data



The copy, serialization and disk I/O makes it slow

In Memory Computing

Use memory to store data



10-100 times faster than the disk approach

Issues

- Is memory large enough to contain all the data?
- How about the cost comparing to disks?
- Fault-tolerance with memory?
- How to represent the data in memory efficiently?

Two traditional abstractions

- Distributed Shared Memory
 - Extend the virtual memory
 - `read/write #address`
- Distributed KV stores
 - `Value = Get(key)`
 - `Put(key, value)`

How to tolerate fault in distributed k-v stores?

A Log Sample

TID	T1	T1	T1	T2	T2	T2	T3
OLD	A=0	B=0		A=100	B=50		A=80
NEW	A=100	B=50	COMMIT	A=80	B=70	COMMIT	A=110

```
begin // T1
A = 100
B = 50
commit // A=100; B=50

begin // T2
A = A - 20
B = B + 20
commit // A=80; B=70

begin // T3
A = A + 30 // A=110
--CRASH--
```

Fault tolerance mechanism of DSM or Key-value stores

- Replica
 - High cost
 - Big time overhead
- Log
 - Big time overhead
- Remote memory access
 - 10-100 times slower than local memory accesses
 - 1-10us vs 100ns

A more lightweight fault tolerance mechanism for in-memory computing is desired

Five Operations Involved

- **1. begin:** allocate a new transaction ID
- **2. write variable:** append an entry to the log
- **3. read variable:** scan the log looking for last committed value
 - As an aside: how to see your own updates?
 - Read uncommitted values from your own TID

Five Operations Involved

Performance Problem of Log-only Approach

- **Write performance** is probably good
 - Sequential writes, instead of random
 - Need to write to disk **twice** instead of **once**
- **Read performance** is **terrible!**
 - Scan the log for every read!
- **Recovery** is instantaneous
 - Nothing to do

- **4. Commit:** write a commit record
 - Expectedly, writing a commit record is the "commit point" for action, because of the way read works (looks for commit record)
 - However, writing log records better be **all-or-nothing**
 - One approach, from last time: make each record fit within one sector
- **5. Abort:** Do nothing or write a record?
 - Could write an abort record, but not strictly needed
 - *Recover from a crash:* do nothing

RDD (Resilient Distributed Datasets)

Solution

RDD (Resilient Distributed Datasets)

- Based on data sets, instead of single data
- Generated with deterministic coarse grained operations (map, filter, join etc.)
- Immutable
- To modify data, create new datasets by transformation

Efficient Fault Tolerance

- Once the data is generated deterministically, and immutable
 - Recover the data from “recompute”
 - Only need to log the sequences to generate rdd. small cost if there's no fault happens

```
messages = textFile(...).filter(_.contains("error"))
           .map(_.split("\t")(2))
```



How to create RDDs

- From Collection


```
val a = sc.parallelize(1 to 9, 3)
```
- from File


```
val b = sc.textFile("/path/to/file")
```
- from other RDDs


```
val c = b.map(line => line.split(","))
```

Spark Operators

Transformations (define a new RDD)	map filter sample groupByKey reduceByKey sortByKey	flatMap union join cogroup cross mapValues
Actions (return a result to driver program)	collect reduce count save lookupKey	

Spark Programming

Scala - Functions

- def helloWord()={ println("Hello World")}
- def addOne(x:Int) = {


```
x+1
```

 }
- def addOne(x:Int) = x+1
- val addOne = (x:Int) => x + 1

RDD operations

- Transformation: generate new RDDs from existing RDDs
 - map, filter, groupBy, sort, distinct, sample...
- Action: return a value from RDD
 - **reduce**, count, collect, first, foreach..., ,

Map

- Map each element in source RDD to ONE element in destination RDD

```
scala> val a = sc.parallelize(1 to 9, 3)
scala> val b = a.map(x => x*2)
scala> a.collect
res10: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9)
scala> b.collect
res11: Array[Int] = Array(2, 4, 6, 8, 10, 12, 14, 16, 18)
reduceByKey
```

- Reduce for key-value list
- Values with the same keys are reduced, and form a new key-value

```
scala> val a = sc.parallelize(List((1,2),(3,4),(3,6)))
scala> a.reduceByKey((x,y) => x + y).collect
res7: Array[(Int, Int)] = Array((1,2), (3,10))
```

Use of _

- val a = (1 to 9)
- Val b = a.map(x => x*2)

is the same as val b = a.map(_ * 2)
- val a = sc.parallelize(List((1,2),(3,4),(3,6))) **a.reduceByKey((x,y) => x + y).collect** is the same as **a.reduceByKey(_ + _).collect**

flatMap

- One element in source rdd will be transformed to multiple elements(0..n) in destination rdd

```
scala> val a = sc.parallelize(1 to 4, 2)
scala> val b = a.flatMap(x => 1 to x)
scala> b.collect
res0: Array[Int] = Array(1, 1, 2, 1, 2, 3, 1, 2, 3, 4)
```

Join

(k,v1) join (k, v2) → (k, (v1, v2))

```
scala> val kv1=sc.parallelize(List(("A",1),("B",2),("C",3)))
scala> val kv3=sc.parallelize(List(("A",4),("A",5),("B",3),("D",30)))
scala> kv1.join(kv3).collect
Array[(String, (Int, Int))] = Array((A,(1,4)), (A,(1,5)), (B,(2,3)))
scala> kv1.union(kv3).collect
Array[(String, (Int, Int))] = Array((A,1), (B,2), (C,3), (A,4), (A,5), (B,3), (D,30))
```

```
scala> val c = sc.parallelize(1 to 10)
scala> c.reduce((x, y) => x + y)
```

res4: Int = 55

Output

```
• saveAsTextFile/saveAsSequenceFile/saveAsObjectFile
```

```
scala> val a = sc.parallelize(List((1,2),(3,4),(3,6)))
scala> a.reduceByKey((x,y) => x + y).saveAsTextFile("test")
```

Join

(k,v1) join (k, v2) → (k, (v1, v2))

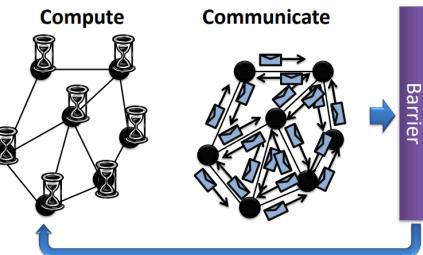
```
scala> val kv1=sc.parallelize(List(("A",1),("B",2),("C",3)))
scala> val kv3=sc.parallelize(List(("A",4),("A",5),("B",3),("D",30)))
scala> kv1.join(kv3).collect
Array[(String, (Int, Int))] = Array((A,(1,4)), (A,(1,5)), (B,(2,3)))
scala> kv1.union(kv3).collect
Array((A,1), (B,2), (C,3), (A,4), (A,5), (B,3), (D,30))
```


Graph Analysis

Spark isn't suited for extreme-scale graph processing

- Slow and memory inefficient

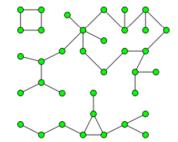
The BSP model does not support asynchronous computing



Weakly Connected Components

Algorithm

- Give each vertex an ID and look at all edges and find the vertex ID who holds the highest value.
- Send the greatest vertex ID to all the edges.
- Process incoming messages and compare it against the value you hold.
- Repeat the process until there are no messages being sent.



What happens if this algorithm is executed asynchronously?

Asynchronous processing

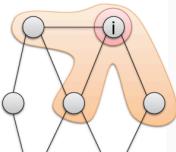
The GraphLab Abstraction

Vertex-Programs directly **read** the neighbors state

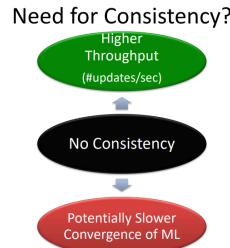
```
GraphLab_PageRank(i)
    // Compute sum over neighbors
    total = 0
    foreach( j in in_neighbors(i)):
        total = total + R[j] * wji

    // Update the PageRank
    R[i] = 0.15 + total

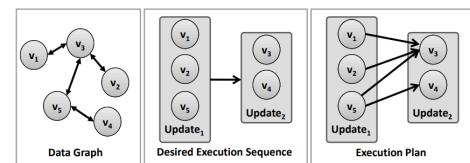
    // Trigger neighbors to run again
    if R[i] not converged then
        foreach( j in out_neighbors(i)):
            signal vertex-program on j
```



Ensuring Race-Free Code in Async processing:

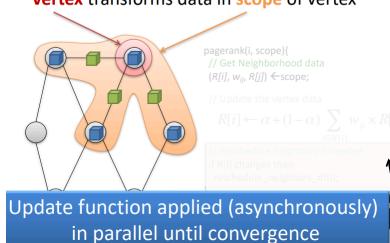


An execution plan example



Update Functions

User-defined program: applied to **vertex** transforms data in **scope** of vertex

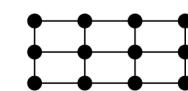


Graph partitioning on power law graphs and GAS model

Needed because Graphs can become unhandable huge.

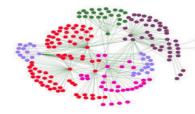
Assumptions of Graph-Parallel Abstractions

Idealized Structure



- **Small neighborhoods**
 - Low degree vertices
- Vertices have similar degree
- Easy to partition

Natural Graph



- **Large Neighborhoods**
 - High degree vertices
- **Power-Law** degree distribution
- **Difficult to partition**

Random Partitioning

- Both GraphLab 1 and Pregel proposed Random (hashed) partitioning for Natural Graphs

$$\text{For } p \text{ Machines: } \mathbb{E} \left[\frac{|\text{Edges Cut}|}{|E|} \right] = 1 - \frac{1}{p}$$

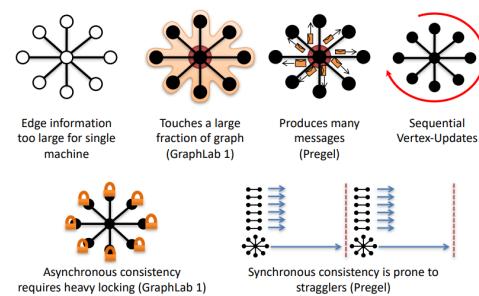
10 Machines → 90% of edges cut
100 Machines → 99% of edges cut!

In Summary

GraphLab 1 and Pregel are not well suited for natural graphs

Problem:

High Degree Vertices Limit Parallelism



- Poor performance on high-degree vertices
- Low Quality Partitioning



Minimizing Communication in GraphLab 2: Vertex Cuts

- Distribute a single vertex-update
 - Move computation to data
 - Parallelize high-degree vertices
- Vertex Partitioning
 - Simple online approach, effectively partitions large power-law graphs

A Common Pattern for Vertex-Programs

```
GraphLab_PageRank(i)
// Compute sum over neighbors
total = 0
foreach( j in in_neighbors(i)):
  total = total + R[j] * wji

// Update the PageRank
R[i] = 0.1 + total

// Trigger neighbors to run again
if R[i] not converged then
  foreach( j in out_neighbors(i))
    signal vertex-program on j
```

Gather Information About Neighborhood

Update Vertex

Signal Neighbors & Modify Edge Data

PageRank in PowerGraph

$$R[i] = 0.15 + \sum_{j \in \text{Nbrs}(i)} w_{ji} R[j]$$

PowerGraph_PageRank(i)

```
Gather(j → i) : return wji * R[j]
sum(a, b) : return a + b;
```

```
Apply(i, Σ) : R[i] = 0.15 + Σ
```

Scatter(i → j) :

if $R[i]$ changed then trigger j to be recomputed

PowerLyra

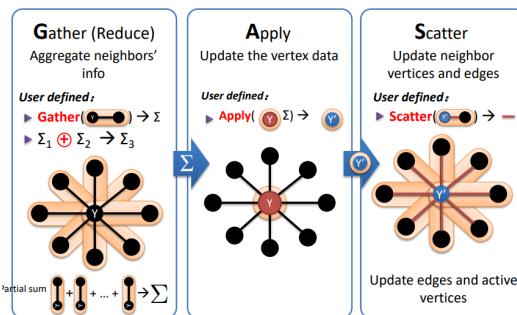
- PowerGraph partitions graph for high degree vertices, but is not friendly to the low degree vertices(locality)
- Optimize the low degree vertices' locality
 - Avoids unnecessary communications

A vertex-cut minimizes
machines per vertex

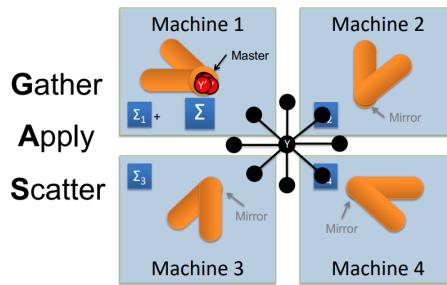
Percolation theory suggests Power Law graphs can be split by removing only a small set of vertices [Albert et al. 2000]

Small vertex cuts possible!

GAS



Distributed GAS



Constructing Vertex-Cuts

- Goal:** Parallel graph partitioning on ingress
- GraphLab 2 provides three simple approaches:
 - Random Edge Placement**
 - Edges are placed randomly by each machine
 - Good theoretical guarantees
 - Greedy Edge Placement with Coordination**
 - Edges are placed using a shared objective
 - Better theoretical guarantees
 - Oblivious-Greedy Edge Placement**
 - Edges are placed using a local objective

Out of core graph processing

The key of out-of-core graph systems is the sequential access to graph data

Generalized Graph Processing

```
for each vertex v
  if v is active
    for each edge e to v
      gather update along e
      apply updated value to v
      Vertex 1.0 1.0 1.0 1.0
```

Edge	1 → 2	1 → 3	2 → 1	2 → 3	3 → 2	4 → 2
	2 → 1	2 → 4	1 → 2	1 → 3	4 → 3	2 → 4
	3 → 2	4 → 3	2 → 3	1 → 3	4 → 3	
	4 → 2	4 → 3	2 → 4			

for each vertex v
 if v has update
 for each edge e from v
 scatter update along e

Can we use just one copy of the edge list?
Do we need all vertices in memory?

Out-of-core Graph Processing

- Processing large graphs with limited memory
 - Low cost graph analysis engine

Possible Solutions

1. Use SSD as a memory extension? (Hadoop, MapReduce)
 - Too many small sequential memory transfers / sec.
2. Compress the graph structure to fit into RAM?
 - Associated values don't compress well, and are mutated.
3. Cluster the graph and handle each cluster separately in RAM?
 - Expensive. The number of inter-cluster edges is big.
4. Caching of hot nodes?
 - Unpredictable performance.

Generalized Graph Processing

```
for each vertex v
  if v is active
    for each edge e to v
      gather update along e
      apply updated value to v
      Vertex 1.0 1.0 1.0 1.0
```

Edge	1 → 2	1 → 3	2 → 1	2 → 3	3 → 2	4 → 2
	2 → 1	2 → 4	1 → 2	1 → 3	4 → 3	2 → 4
	3 → 2	4 → 3	2 → 3	1 → 3	4 → 3	
	4 → 2	4 → 3	2 → 4			

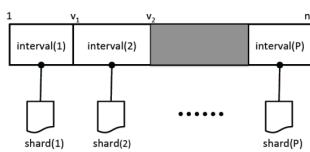
for each vertex v
 if v has update
 for each edge e from v
 scatter update along e

Gather&Apply can be separated from Scatter

GRAPHCHI – Parallel Sliding Window

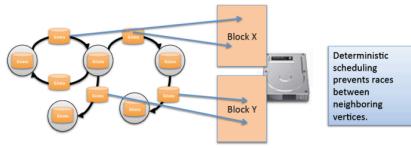
PSW: Shards and Intervals

- Vertices are numbered from 1 to n
 - P intervals, each associated with a shard on disk.
 - sub-graph = interval of vertices



PSW: Execute updates

- Update-function is executed on interval's vertices
- Edges have pointers to the loaded data blocks
- Changes take effect immediately → asynchronous.



Graphchi processing model

for each vertex partition (Intervals)
load all vertex and in/out edges

for each vertex v
if v is active
for each edge e to v
gather update along e
apply updated value to v

for each vertex v
if v has update
for each edge e from v
scatter update along e by attaching the info

Write all in/out edges

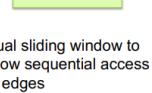
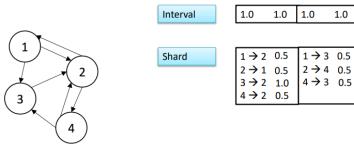
Example



How is gather and scatter implemented?

Example with (simplified)

PageRank



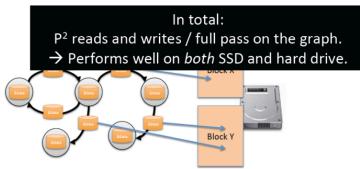
Dual sliding window to allow sequential access to edges

Remove the need to sort in preprocessing

2-D partition to allow selective scheduling on edges; In place apply to avoid shuffle

PSW: Commit to Disk

- In write phase, the blocks are written back to disk
 - Next load-phase sees the preceding writes → asynchronous.



GraphChi is good

- But have two major drawbacks

- P^2 I/O are required
 - Poor performance when P is large

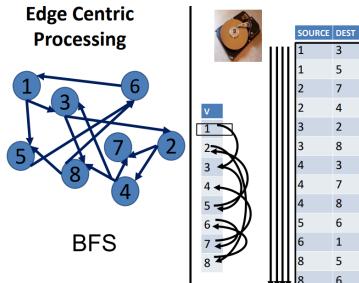
- High preprocessing time because sorting is required for each shard

XSTREAM – Edge Centric

X-Stream : Edge-Centric model

for each vertex v
if v has update
for each edge e from v
scatter update along e

for each edge e
If e.src has update
scatter update along e

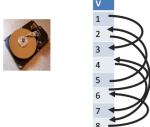


Vertex-Centric → Edge-Centric

Streaming Partitions

Streaming Partitions

- Problem: still have random access to vertex set



- Solution: partition the graph into streaming partitions

Streaming Partitions

- A streaming partition is
 - A subset of the vertices that fits in RAM
 - All edges whose source vertex is in that subset
 - No requirement on quality of the partition

The issues of X-Stream

- X-Stream
 - Separate gather/scatter phase, with shuffle between them, introduce many extra I/O operations
 - Poor performance for BFS and WCC-like algorithms
 - Only a small portion of edges are touched for some iterations, but they would scan all edges

Tradeoff

Vertex-centric Scatter-Gather: Random Access Bandwidth

- Edge-centric Scatter-Gather
 - #of scatters * Edge data / sequential access bandwidth
- Sequential Access Bandwidth >> Random Access Bandwidth
- Few scatter-gather iterations for real world graphs
 - Give an example that is not suitable for this system?

Edge Data/

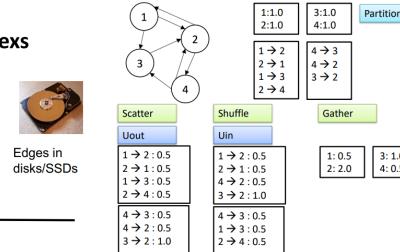
Partition Edges and vertices

V1	SOURCE	DEST
1	5	
2	7	
3	2	
4	8	
4	3	
4	7	
4	8	
5	6	
5	1	
6	1	
6	8	
7	2	
7	4	
8	5	
8	6	

Edges in disks/SSDs

V2	SOURCE	DEST
5	6	
6	8	
7	5	
8	6	

Example on X-Stream



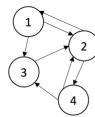
GRIDGRAPH – 2D Partition

GridGraph

- A 2-level partitioning graph data structure
 - Improve locality
 - Enable more effective streaming on both edges and vertices
- A Streaming-Apply model
 - Combines the gather-apply-scatter operations
- A programming abstraction
 - To enable the streaming-apply model

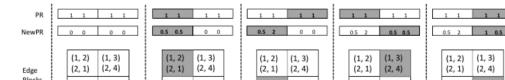
The Grid Representation

- 2-level hierarchical partitioning
 - 1st dimension
 - partition the vertices into chunks(1,2)(3,4)
 - partition the edges into shards by source vertex((1,2),(2,1),(1,3),(2,4))
 - 2nd dimension
 - partition the shards into blocks by destination vertex



(1, 2)	(1, 3)
(2, 1)	(2, 4)
(3, 2)	(4, 3)
(4, 2)	

The Streaming-Apply Model



Dual Sliding Windows:

- 1 read window + 1 write window
- Contents of active windows can be fit into RAM

Each iteration, just read edge list once, has opportunity to skip some locks

```
Algorithm 1 Edge Block Streaming
for j ← 1, P do
  for i ← 1, P do
    if ChunksActive(i) then
      StreamEdgeBlock(i,j);
    end if
  end for
end for
```

Streaming-Apply Processing Model

- Vertex accesses are aggregated
 - Good locality
 - Only 2 partitions of vertices are accessed within each edge block
- On-the-fly updates onto vertices
 - Reduces I/O
 - Enables asynchronous implementation of algorithms (like WCC, etc.) which converges faster

GridGraph processing model

- for each vertex partition
 - load all vertex and in/out edges
- for each edge e
 - gather update along e
 - read e.src
 - apply updated value to e.dst

Write the unnecessary vertex partition back to disk

Programming Abstraction

Algorithm 2 Edge Streaming Interface

```
function STREAMEDGES(Fe, F)
  Sum = 0
  for all block do
    for all edge ∈ block do
      if F(edge.source) then
        Sum += Fe(edge)
      end if
    end for
  end for
  return Sum
end function
```



StreamEdges(Fe, F):
Fe: process the edge, which is allowed to side-effect related vertex data
F: check whether the edge should be processed (by looking at the source vertex)

Algorithm 3 Vertex Streaming Interface

```
function STREAMVERTICES(Fv, F)
  Sum = 0
  for all vertex do
    if F(vertex) then
      Sum += Fv(vertex)
    end if
  end for
  return Sum
end function
```



StreamVertices(Fv, F):
Fv: process the vertex
F: check if the vertex should be processed

BFS

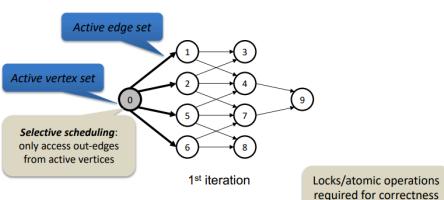
```
Algorithm 4 BFS
function ISACTIVE(v)
  return ActiveIn[v]
end function
function VISIT(e)
  if CAS(&Parent[e.dest], -1, e.source) then
    ActiveOut[e.dest] = 1
    return 1
  end if
  return 0
end function
ActiveVertices = 1
Parent = {-1, ..., -1}
ActiveIn = {0, ..., 0}
Parent[start] = start
ActiveIn[start] = 1
while ActiveVertices > 0 do
  ActiveOut = {0, ..., 0}
  ActiveVertices = StreamEdges(VISIT, ISACTIVE)
  Swap(ActiveIn, ActiveOut)
end while
```

PageRank Implementation

```
Algorithm 3 PageRank
function CONTRIBUTE(e)
  Accum(&NewPR[e.dest], PR[e.source])
end function
function COMPUTE(v)
  NewPR[v] = 1 - d + d × NewPR[v]
  return |NewPR[v] - PR[v]|
end function
d = 0.85
PR = {1, ..., 1}
Converged = 0
while ¬Converged do
  NewPR = {0, ..., 0}
  StreamEdges(Contribute)
  Diff = StreamVertices(Compute)
  Swap(PR, NewPR)
  Converged = Diff / V ≤ Threshold
end while
```

Top Down vs. Bottom Up

BFS Example (top-down model)



function top-down-step(vertices, frontier, next, parents)

```
for v ∈ frontier do
  for n ∈ neighbors[v] do
    if parents[n] = -1 then
      parents[n] ← v
      next ← next ∪ {n}
    end if
  end for
end for
```

Heuristics for mode switch

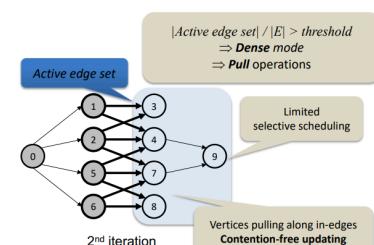
Dual mode updates proposed in shared-memory systems (LigrapPP+PP+13)

```
|Active edge set| / |E| < threshold
⇒ Sparse mode
⇒ Push operations
```

```
|Active edge set| / |E| > threshold
⇒ Dense mode
⇒ Pull operations
```

- Top-down
 - Less active edges – sparse – push
- Bottom-up
 - More active edges – dense – pull

BFS(Bottom-up mode)



Bottom-up BFS

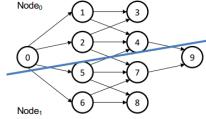
```
function bottom-up-step(vertices, frontier, next, parents)
  for v ∈ vertices do
    if parents[v] = -1 then
      for n ∈ neighbors[v] do
        if n ∈ frontier then
          parents[v] ← n
          next ← next ∪ {v}
        break
      end for
    end if
  end for
```

GEMINI – Distributed push/pull

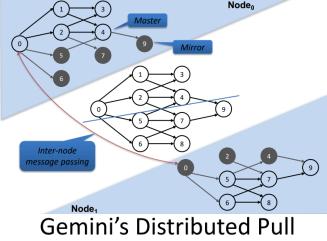
We Propose: *Gemini*

- Build **scalability** on top of **efficiency**
 - Avoid unnecessary “distributed” side-effects
 - Optimize computation on partitioned graphs
- Shift of design focus
 - Designed for distributed, but **computation-centric**
 - Modern clusters have fast interconnects
 - Computation-communication overlap in place
 - Efficiency optimizations
 - Adaptive push-/pull-style computation
 - Hierarchical chunk-based partitioning
 - Scalability optimizations
 - Locality-aware chunking
 - Chunk-based work-stealing

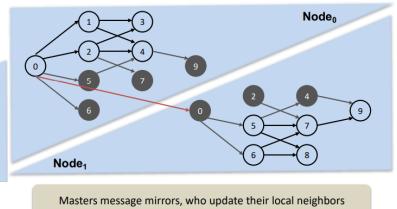
Distributed Dual-mode Computation



When Distributed to 2 nodes



Gemini’s Distributed Push



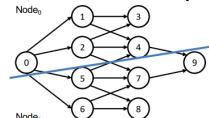
More Benefits of Chunking

- No need to convert vertex IDs (global \leftrightarrow local)
 - Simple bookkeeping of partition information
 - $O(p)$ chunk boundaries
 - Simple management of vertex data arrays
 - Just allocate entire array in shared memory



- Can be applied recursively at different levels

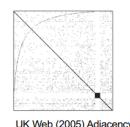
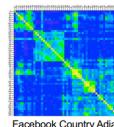
Gemini’s Choice for Graph Partition



- Chunks
 - Divide vertex set V into p contiguous chunks
 - Dual-mode edge data distributed accordingly

Why Chunk-based Partitioning?

- It preserves locality!
 - Fact: **locality** exists in many real-world graphs
 - E.g., WebGraph[WWW’04], BLP[WSOM’13]
 - Vertices “semantically” ordered
 - Small distance between source and destination vertices
- Preprocessing affordable when vertices unordered
 - E.g., BFS[Algorithms’09], LLPI[WWW’11]

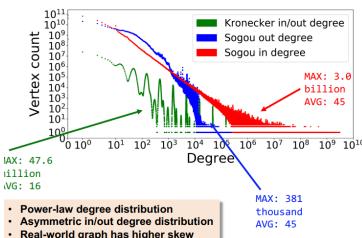


Locality-aware Chunking

- Where to draw partition lines?
 - Balancing by edges?
 - Chunk size affects random access efficiency!
- Gemini considers both vertex and edge
 - Edge: the amount of work to be processed
 - Vertex: the processing speed of work (locality)
 - Hybrid metric: $\alpha \cdot |V_i| + |E_i|$
 - α set to 8/(p-1) currently default value

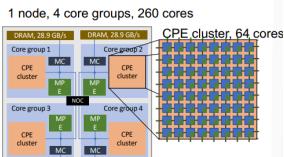
Shentu - Designed for Extreme Scale Graph

Problem 1: Severe Load Imbalance



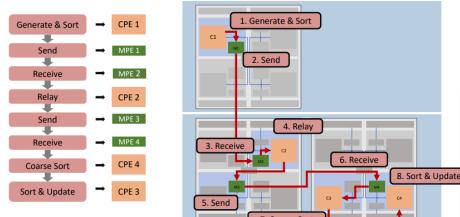
Problem 3: Heterogeneous Architecture

- MPE (Management Processing Element)
 - General-purpose processor
 - Low memory bandwidth
- CPE (Compute Processing Element)
 - High memory bandwidth
 - Limited functionality
- Shared memory between MPE and CPE

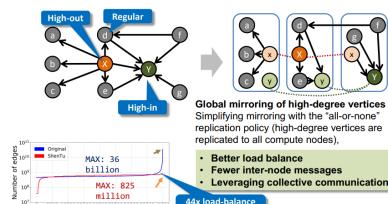


Requires careful task mapping and coordination

Solution 3a: Heterogeneous Task Mapping

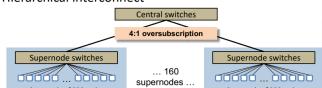


Solution 1: Degree Aware Messaging



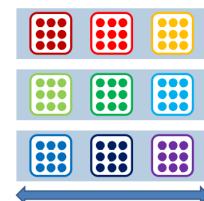
Problem 2a: Too Many Small Messages

- Traditional distributed graph computing model (1D)
 - $O(N \times N)$ messages
 - Not suitable to extreme-scale processing
 - Too expensive on 40,000+ nodes
 - # nodes >> average degree: little opportunity for combining messages
- Hierarchical interconnect

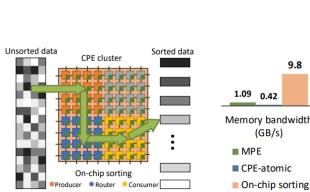


Solution 2a: Supernode Routing

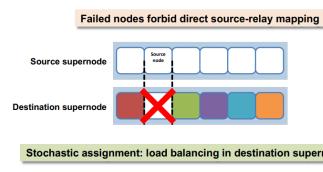
- Stage 1
 - Messages grouped by destination supernodes
 - Combined messages sent to relay node in destination supernode
- Stage 2
 - Messages grouped by destination nodes within same supernode
 - Combined messages sent to final destination node



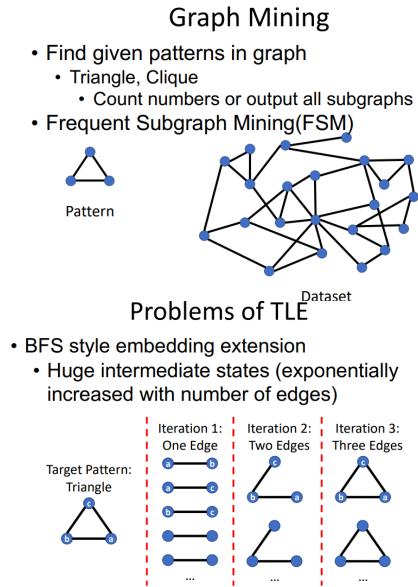
Solution 3b: CPE Cluster On-chip Sorting



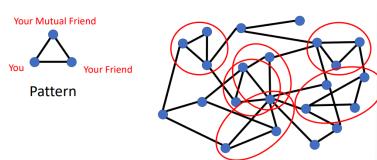
Solution 2b: Topology-aware Supernode Relaying



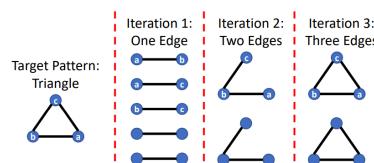
Graph Mining



Triangle Counting



Think Like a Embedding(TLE)



Triangle Counting is well studied

You Mutual Friend
You Your Friend
Pattern

```
Algorithm 1: Triangle counting.
input : G : the Graph.
output : n : the number of triangles in G.
1 begin
2   n ← 0;
3   for  $v_i$  in  $V$  do
4     //  $N(v_i)$  returns a set that contains
5     // all  $v_j$ 's neighbors
6     for  $v_i$  in  $N(v_i)$  do
7        $s \leftarrow N(v_i) \cap N(v_j)$ ;
8       n ← n + | $s$ |;
9   n ← n/6;
```

Algorithm 1: Exploration step in Arabesque.

```
input : Set  $I$  of initial embeddings for this step
output: Set  $F$  of extended embeddings for the next step (initially empty)
output: Set  $O$  of new outputs (initially empty)
foreach  $e \in I$  such that  $\alpha(e)$  do
  add  $\beta(e)$  to  $O$ ;
   $C \leftarrow$  set of all extensions of  $e$  obtained by adding one incident
  edge / neighboring vertex;
  foreach  $e' \in C$  do
    if  $\phi(e')$  and there exists no  $e'' \in F$  automorphic to  $e'$  then
      add  $\pi(e')$  to  $O$ ;
      add  $e'$  to  $F$ ;
run aggregation functions;
```

Stream Processing

What we did so far: Batch processing, processing all data

- Batch processing, processing all data
 - Map-Reduce
 - Spark
 - GraphLab / Gemini



Stream Processing

- $F(X + \Delta X) = F(X) \text{ op } H(\Delta X)$
- When we calculate $F(X + \Delta X)$, we don't need to use $X + \Delta X$, instead we only need $F(X)$ and ΔX
- We call this computing paradigm Stream Processing

Real-time / Scalability

- Batch processing
 - Fixed problem size
 - Latency: Tens of hours processing time
- Stream Processing
 - Arrival rate varies heavily
 - Latency: Seconds or even less
 - Processing all coming data or
 - Pre-defined degrading approach

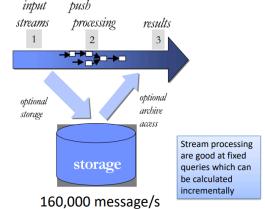
Fault Tolerance

- Batch
 - ETL can remove faults in data
 - Checkpoint and/or recalculation
- Streaming
 - Has to deal with data fault in real time
 - Fault tolerance mechanism must be low-overhead, and maintain the real time demand

Use traditional database



Stream Processing



Queue

- Matching Producers and Consumers
 - Speed matching
 - Interfaces
 - Pub – Subscribe model
 - Persistency
 - Could replay data
 - Performance optimization
 - Batching data (Similar to relay in Shentu)

Naive Implementation

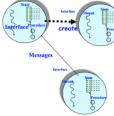
Mismatched Speed of Producer/Consumer : Add a Queue



Input Data speed → Processing Speed

Queue is used to address the TEMPORAL speed mismatch

Actor run asynchronously*



1. message
2. process message
3. Actor 1 sends message to Actor 2
4. sleep
5. Actor 2 receives message
6. process message
7. Actor 2 sends message to Actor 1
8. sleep
9. Actor 1 receives message
10. process message
11. Actor 1 sends message to Actor 2
12. sleep
13. Actor 2 receives message
14. process message
15. Actor 2 sends message to Actor 1
16. sleep
17. Actor 1 receives message
18. process message
19. Actor 1 sends message to Actor 2
20. sleep
21. Actor 2 receives message
22. process message
23. Actor 2 sends message to Actor 1
24. sleep
25. Actor 1 receives message
26. process message
27. Actor 1 sends message to Actor 2
28. sleep
29. Actor 2 receives message
30. process message
31. Actor 2 sends message to Actor 1
32. sleep
33. Actor 1 receives message
34. process message
35. Actor 1 sends message to Actor 2
36. sleep
37. Actor 2 receives message
38. process message
39. Actor 2 sends message to Actor 1
40. sleep
41. Actor 1 receives message
42. process message
43. Actor 1 sends message to Actor 2
44. sleep
45. Actor 2 receives message
46. process message
47. Actor 2 sends message to Actor 1
48. sleep
49. Actor 1 receives message
50. process message
51. Actor 1 sends message to Actor 2
52. sleep
53. Actor 2 receives message
54. process message
55. Actor 2 sends message to Actor 1
56. sleep
57. Actor 1 receives message
58. process message
59. Actor 1 sends message to Actor 2
60. sleep
61. Actor 2 receives message
62. process message
63. Actor 2 sends message to Actor 1
64. sleep
65. Actor 1 receives message
66. process message
67. Actor 1 sends message to Actor 2
68. sleep
69. Actor 2 receives message
70. process message
71. Actor 2 sends message to Actor 1
72. sleep
73. Actor 1 receives message
74. process message
75. Actor 1 sends message to Actor 2
76. sleep
77. Actor 2 receives message
78. process message
79. Actor 2 sends message to Actor 1
80. sleep
81. Actor 1 receives message
82. process message
83. Actor 1 sends message to Actor 2
84. sleep
85. Actor 2 receives message
86. process message
87. Actor 2 sends message to Actor 1
88. sleep
89. Actor 1 receives message
90. process message
91. Actor 1 sends message to Actor 2
92. sleep
93. Actor 2 receives message
94. process message
95. Actor 2 sends message to Actor 1
96. sleep
97. Actor 1 receives message
98. process message
99. Actor 1 sends message to Actor 2
100. sleep
101. Actor 2 receives message
102. process message
103. Actor 2 sends message to Actor 1
104. sleep
105. Actor 1 receives message
106. process message
107. Actor 1 sends message to Actor 2
108. sleep
109. Actor 2 receives message
110. process message
111. Actor 2 sends message to Actor 1
112. sleep
113. Actor 1 receives message
114. process message
115. Actor 1 sends message to Actor 2
116. sleep
117. Actor 2 receives message
118. process message
119. Actor 2 sends message to Actor 1
120. sleep
121. Actor 1 receives message
122. process message
123. Actor 1 sends message to Actor 2
124. sleep
125. Actor 2 receives message
126. process message
127. Actor 2 sends message to Actor 1
128. sleep
129. Actor 1 receives message
130. process message
131. Actor 1 sends message to Actor 2
132. sleep
133. Actor 2 receives message
134. process message
135. Actor 2 sends message to Actor 1
136. sleep
137. Actor 1 receives message
138. process message
139. Actor 1 sends message to Actor 2
140. sleep
141. Actor 2 receives message
142. process message
143. Actor 2 sends message to Actor 1
144. sleep
145. Actor 1 receives message
146. process message
147. Actor 1 sends message to Actor 2
148. sleep
149. Actor 2 receives message
150. process message
151. Actor 2 sends message to Actor 1
152. sleep
153. Actor 1 receives message
154. process message
155. Actor 1 sends message to Actor 2
156. sleep
157. Actor 2 receives message
158. process message
159. Actor 2 sends message to Actor 1
160. sleep
161. Actor 1 receives message
162. process message
163. Actor 1 sends message to Actor 2
164. sleep
165. Actor 2 receives message
166. process message
167. Actor 2 sends message to Actor 1
168. sleep
169. Actor 1 receives message
170. process message
171. Actor 1 sends message to Actor 2
172. sleep
173. Actor 2 receives message
174. process message
175. Actor 2 sends message to Actor 1
176. sleep
177. Actor 1 receives message
178. process message
179. Actor 1 sends message to Actor 2
180. sleep
181. Actor 2 receives message
182. process message
183. Actor 2 sends message to Actor 1
184. sleep
185. Actor 1 receives message
186. process message
187. Actor 1 sends message to Actor 2
188. sleep
189. Actor 2 receives message
190. process message
191. Actor 2 sends message to Actor 1
192. sleep
193. Actor 1 receives message
194. process message
195. Actor 1 sends message to Actor 2
196. sleep
197. Actor 2 receives message
198. process message
199. Actor 2 sends message to Actor 1
200. sleep
201. Actor 1 receives message
202. process message
203. Actor 1 sends message to Actor 2
204. sleep
205. Actor 2 receives message
206. process message
207. Actor 2 sends message to Actor 1
208. sleep
209. Actor 1 receives message
210. process message
211. Actor 1 sends message to Actor 2
212. sleep
213. Actor 2 receives message
214. process message
215. Actor 2 sends message to Actor 1
216. sleep
217. Actor 1 receives message
218. process message
219. Actor 1 sends message to Actor 2
220. sleep
221. Actor 2 receives message
222. process message
223. Actor 2 sends message to Actor 1
224. sleep
225. Actor 1 receives message
226. process message
227. Actor 1 sends message to Actor 2
228. sleep
229. Actor 2 receives message
230. process message
231. Actor 2 sends message to Actor 1
232. sleep
233. Actor 1 receives message
234. process message
235. Actor 1 sends message to Actor 2
236. sleep
237. Actor 2 receives message
238. process message
239. Actor 2 sends message to Actor 1
240. sleep
241. Actor 1 receives message
242. process message
243. Actor 1 sends message to Actor 2
244. sleep
245. Actor 2 receives message
246. process message
247. Actor 2 sends message to Actor 1
248. sleep
249. Actor 1 receives message
250. process message
251. Actor 1 sends message to Actor 2
252. sleep
253. Actor 2 receives message
254. process message
255. Actor 2 sends message to Actor 1
256. sleep
257. Actor 1 receives message
258. process message
259. Actor 1 sends message to Actor 2
260. sleep
261. Actor 2 receives message
262. process message
263. Actor 2 sends message to Actor 1
264. sleep
265. Actor 1 receives message
266. process message
267. Actor 1 sends message to Actor 2
268. sleep
269. Actor 2 receives message
270. process message
271. Actor 2 sends message to Actor 1
272. sleep
273. Actor 1 receives message
274. process message
275. Actor 1 sends message to Actor 2
276. sleep
277. Actor 2 receives message
278. process message
279. Actor 2 sends message to Actor 1
280. sleep
281. Actor 1 receives message
282. process message
283. Actor 1 sends message to Actor 2
284. sleep
285. Actor 2 receives message
286. process message
287. Actor 2 sends message to Actor 1
288. sleep
289. Actor 1 receives message
290. process message
291. Actor 1 sends message to Actor 2
292. sleep
293. Actor 2 receives message
294. process message
295. Actor 2 sends message to Actor 1
296. sleep
297. Actor 1 receives message
298. process message
299. Actor 1 sends message to Actor 2
300. sleep
301. Actor 2 receives message
302. process message
303. Actor 2 sends message to Actor 1
304. sleep
305. Actor 1 receives message
306. process message
307. Actor 1 sends message to Actor 2
308. sleep
309. Actor 2 receives message
310. process message
311. Actor 2 sends message to Actor 1
312. sleep
313. Actor 1 receives message
314. process message
315. Actor 1 sends message to Actor 2
316. sleep
317. Actor 2 receives message
318. process message
319. Actor 2 sends message to Actor 1
320. sleep
321. Actor 1 receives message
322. process message
323. Actor 1 sends message to Actor 2
324. sleep
325. Actor 2 receives message
326. process message
327. Actor 2 sends message to Actor 1
328. sleep
329. Actor 1 receives message
330. process message
331. Actor 1 sends message to Actor 2
332. sleep
333. Actor 2 receives message
334. process message
335. Actor 2 sends message to Actor 1
336. sleep
337. Actor 1 receives message
338. process message
339. Actor 1 sends message to Actor 2
340. sleep
341. Actor 2 receives message
342. process message
343. Actor 2 sends message to Actor 1
344. sleep
345. Actor 1 receives message
346. process message
347. Actor 1 sends message to Actor 2
348. sleep
349. Actor 2 receives message
350. process message
351. Actor 2 sends message to Actor 1
352. sleep
353. Actor 1 receives message
354. process message
355. Actor 1 sends message to Actor 2
356. sleep
357. Actor 2 receives message
358. process message
359. Actor 2 sends message to Actor 1
360. sleep
361. Actor 1 receives message
362. process message
363. Actor 1 sends message to Actor 2
364. sleep
365. Actor 2 receives message
366. process message
367. Actor 2 sends message to Actor 1
368. sleep
369. Actor 1 receives message
370. process message
371. Actor 1 sends message to Actor 2
372. sleep
373. Actor 2 receives message
374. process message
375. Actor 2 sends message to Actor 1
376. sleep
377. Actor 1 receives message
378. process message
379. Actor 1 sends message to Actor 2
380. sleep
381. Actor 2 receives message
382. process message
383. Actor 2 sends message to Actor 1
384. sleep
385. Actor 1 receives message
386. process message
387. Actor 1 sends message to Actor 2
388. sleep
389. Actor 2 receives message
390. process message
391. Actor 2 sends message to Actor 1
392. sleep
393. Actor 1 receives message
394. process message
395. Actor 1 sends message to Actor 2
396. sleep
397. Actor 2 receives message
398. process message
399. Actor 2 sends message to Actor 1
400. sleep
401. Actor 1 receives message
402. process message
403. Actor 1 sends message to Actor 2
404. sleep
405. Actor 2 receives message
406. process message
407. Actor 2 sends message to Actor 1
408. sleep
409. Actor 1 receives message
410. process message
411. Actor 1 sends message to Actor 2
412. sleep
413. Actor 2 receives message
414. process message
415. Actor 2 sends message to Actor 1
416. sleep
417. Actor 1 receives message
418. process message
419. Actor 1 sends message to Actor 2
420. sleep
421. Actor 2 receives message
422. process message
423. Actor 2 sends message to Actor 1
424. sleep
425. Actor 1 receives message
426. process message
427. Actor 1 sends message to Actor 2
428. sleep
429. Actor 2 receives message
430. process message
431. Actor 2 sends message to Actor 1
432. sleep
433. Actor 1 receives message
434. process message
435. Actor 1 sends message to Actor 2
436. sleep
437. Actor 2 receives message
438. process message
439. Actor 2 sends message to Actor 1
440. sleep
441. Actor 1 receives message
442. process message
443. Actor 1 sends message to Actor 2
444. sleep
445. Actor 2 receives message
446. process message
447. Actor 2 sends message to Actor 1
448. sleep
449. Actor 1 receives message
450. process message
451. Actor 1 sends message to Actor 2
452. sleep
453. Actor 2 receives message
454. process message
455. Actor 2 sends message to Actor 1
456. sleep
457. Actor 1 receives message
458. process message
459. Actor 1 sends message to Actor 2
460. sleep
461. Actor 2 receives message
462. process message
463. Actor 2 sends message to Actor 1
464. sleep
465. Actor 1 receives message
466. process message
467. Actor 1 sends message to Actor 2
468. sleep
469. Actor 2 receives message
470. process message
471. Actor 2 sends message to Actor 1
472. sleep
473. Actor 1 receives message
474. process message
475. Actor 1 sends message to Actor 2
476. sleep
477. Actor 2 receives message
478. process message
479. Actor 2 sends message to Actor 1
480. sleep
481. Actor 1 receives message
482. process message
483. Actor 1 sends message to Actor 2
484. sleep
485. Actor 2 receives message
486. process message
487. Actor 2 sends message to Actor 1
488. sleep
489. Actor 1 receives message
490. process message
491. Actor 1 sends message to Actor 2
492. sleep
493. Actor 2 receives message
494. process message
495. Actor 2 sends message to Actor 1
496. sleep
497. Actor 1 receives message
498. process message
499. Actor 1 sends message to Actor 2
500. sleep
501. Actor 2 receives message
502. process message
503. Actor 2 sends message to Actor 1
504. sleep
505. Actor 1 receives message
506. process message
507. Actor 1 sends message to Actor 2
508. sleep
509. Actor 2 receives message
510. process message
511. Actor 2 sends message to Actor 1
512. sleep
513. Actor 1 receives message
514. process message
515. Actor 1 sends message to Actor 2
516. sleep
517. Actor 2 receives message
518. process message
519. Actor 2 sends message to Actor 1
520. sleep
521. Actor 1 receives message
522. process message
523. Actor 1 sends message to Actor 2
524. sleep
525. Actor 2 receives message
526. process message
527. Actor 2 sends message to Actor 1
528. sleep
529. Actor 1 receives message
530. process message
531. Actor 1 sends message to Actor 2
532. sleep
533. Actor 2 receives message
534. process message
535. Actor 2 sends message to Actor 1
536. sleep
537. Actor 1 receives message
538. process message
539. Actor 1 sends message to Actor 2
540. sleep
541. Actor 2 receives message
542. process message
543. Actor 2 sends message to Actor 1
544. sleep
545. Actor 1 receives message
546. process message
547. Actor 1 sends message to Actor 2
548. sleep
549. Actor 2 receives message
550. process message
551. Actor 2 sends message to Actor 1
552. sleep
553. Actor 1 receives message
554. process message
555. Actor 1 sends message to Actor 2
556. sleep
557. Actor 2 receives message
558. process message
559. Actor 2 sends message to Actor 1
560. sleep
561. Actor 1 receives message
562. process message
563. Actor 1 sends message to Actor 2
564. sleep
565. Actor 2 receives message
566. process message
567. Actor 2 sends message to Actor 1
568. sleep
569. Actor 1 receives message
570. process message
571. Actor 1 sends message to Actor 2
572. sleep
573. Actor 2 receives message
574. process message
575. Actor 2 sends message to Actor 1
576. sleep
577. Actor 1 receives message
578. process message
579. Actor 1 sends message to Actor 2
580. sleep
581. Actor 2 receives message
582. process message
583. Actor 2 sends message to Actor 1
584. sleep
585. Actor 1 receives message
586. process message
587. Actor 1 sends message to Actor 2
588. sleep
589. Actor 2 receives message
590. process message
591. Actor 2 sends message to Actor 1
592. sleep
593. Actor 1 receives message
594. process message
595. Actor 1 sends message to Actor 2
596. sleep
597. Actor 2 receives message
598. process message
599. Actor 2 sends message to Actor 1
600. sleep
601. Actor 1 receives message
602. process message
603. Actor 1 sends message to Actor 2
604. sleep
605. Actor 2 receives message
606. process message
607. Actor 2 sends message to Actor 1
608. sleep
609. Actor 1 receives message
610. process message
611. Actor 1 sends message to Actor 2
612. sleep
613. Actor 2 receives message
614. process message
615. Actor 2 sends message to Actor 1
616. sleep
617. Actor 1 receives message
618. process message
619. Actor 1 sends message to Actor 2
620. sleep
621. Actor 2 receives message
622. process message
623. Actor 2 sends message to Actor 1
624. sleep
625. Actor 1 receives message
626. process message
627. Actor 1 sends message to Actor 2
628. sleep
629. Actor 2 receives message
630. process message
631. Actor 2 sends message to Actor 1
632. sleep
633. Actor 1 receives message
634. process message
635. Actor 1 sends message to Actor 2
636. sleep
637. Actor 2 receives message
638. process message
639. Actor 2 sends message to Actor 1
640. sleep
641. Actor 1 receives message
642. process message
643. Actor 1 sends message to Actor 2
644. sleep
645. Actor 2 receives message
646. process message
647. Actor 2 sends message to Actor 1
648. sleep
649. Actor 1 receives message
650. process message
651. Actor 1 sends message to Actor 2
652. sleep
653. Actor 2 receives message
654. process message
655. Actor 2 sends message to Actor 1
656. sleep
657. Actor 1 receives message
658. process message
659. Actor 1 sends message to Actor 2
660. sleep
661. Actor 2 receives message
662. process message
663. Actor 2 sends message to Actor 1
664. sleep
665. Actor 1 receives message
666. process message
667. Actor 1 sends message to Actor 2
668. sleep
669. Actor 2 receives message
670. process message
671. Actor 2 sends message to Actor 1
672. sleep
673. Actor 1 receives message
674. process message
675. Actor 1 sends message to Actor 2
676. sleep
677. Actor 2 receives message
678. process message
679. Actor 2 sends message to Actor 1
680. sleep
681. Actor 1 receives message
682. process message
683. Actor 1 sends message to Actor 2
684. sleep
685. Actor 2 receives message
686. process message
687. Actor 2 sends message to Actor 1
688. sleep
689. Actor 1 receives message
690. process message
691. Actor 1 sends message to Actor 2
692. sleep
693. Actor 2 receives message
694. process message
695. Actor 2 sends message to Actor 1
696. sleep
697. Actor 1 receives message
698. process message
699. Actor 1 sends message to Actor 2
700. sleep
701. Actor 2 receives message
702. process message
703. Actor 2 sends message to Actor 1
704. sleep
705. Actor 1 receives message
706. process message
707. Actor 1 sends message to Actor 2
708. sleep
709. Actor 2 receives message
710. process message
711. Actor 2 sends message to Actor 1
712. sleep
713. Actor 1 receives message
714. process message
715. Actor 1 sends message to Actor 2
716. sleep
717. Actor 2 receives message
718. process message
719. Actor 2 sends message to Actor 1
720. sleep
721. Actor 1 receives message
722. process message
723. Actor 1 sends message to Actor 2
724. sleep
725. Actor 2 receives message
726. process message
727. Actor 2 sends message to Actor 1
728. sleep
729. Actor 1 receives message
730. process message
731. Actor 1 sends message to Actor 2
732. sleep
733. Actor 2 receives message
734. process message
735. Actor 2 sends message to Actor 1
736. sleep
737. Actor 1 receives message
738. process message
739. Actor 1 sends message to Actor 2
740. sleep
741. Actor 2 receives message
742. process message
743. Actor 2 sends message to Actor 1
744. sleep
745. Actor 1 receives message
746. process message
747. Actor 1 sends message to Actor 2
748. sleep
749. Actor 2 receives message
750. process message
751. Actor 2 sends message to Actor 1
752. sleep
753. Actor 1 receives message
754. process message
755. Actor 1 sends message to Actor 2
756. sleep
757. Actor 2 receives message
758. process message
759. Actor 2 sends message to Actor 1
760. sleep
761. Actor 1 receives message
762. process message
763. Actor 1 sends message to Actor 2
764. sleep
765. Actor 2 receives message
766. process message
767. Actor 2 sends message to Actor 1
768. sleep
769. Actor 1 receives message
770. process message
771. Actor 1 sends message to Actor 2
772. sleep
773. Actor 2 receives message
774. process message
775. Actor 2 sends message to Actor 1
776. sleep
777. Actor 1 receives message
778. process message
779. Actor 1 sends message to Actor 2
780. sleep
781. Actor 2 receives message
782. process message
783. Actor 2 sends message to Actor 1
784. sleep
785. Actor 1 receives message
786. process message
787. Actor 1 sends message to Actor 2
788. sleep
789. Actor 2 receives message
790. process message
791. Actor 2 sends message to Actor 1
792. sleep
793. Actor 1 receives message
794. process message
795. Actor 1 sends message to Actor 2
796. sleep
797. Actor 2 receives message
798. process message
799. Actor 2 sends message to Actor 1
800. sleep
801. Actor 1 receives message
802. process message
803. Actor 1 sends message to Actor 2
804. sleep
805. Actor 2 receives message
806. process message
807. Actor 2 sends message to Actor 1
808. sleep
809. Actor 1 receives message
810. process message
811. Actor 1 sends message to Actor 2
812. sleep
813. Actor 2 receives message
814. process message
815. Actor 2 sends message to Actor 1
816. sleep
817. Actor 1 receives message
818. process message
819. Actor 1 sends message to Actor 2
820. sleep
821. Actor 2 receives message
822. process message
823. Actor 2 sends message to Actor 1
824. sleep
825. Actor 1 receives message
826. process message
827. Actor 1 sends message to Actor 2
828. sleep
829. Actor 2 receives message
830. process message
831. Actor 2 sends message to Actor 1
832. sleep
833. Actor 1 receives message
834. process message
835. Actor 1 sends message to Actor 2
836. sleep
837. Actor 2 receives message
838. process message
839. Actor 2 sends message to Actor 1
840. sleep
841. Actor 1 receives message
842. process message
843. Actor 1 sends message to Actor 2
844. sleep
845. Actor 2 receives message
846. process message
847. Actor 2 sends message to Actor 1
848. sleep
849. Actor 1 receives message
850. process message
851. Actor 1 sends message to Actor 2
852. sleep
853. Actor 2 receives message
854. process message
855. Actor 2 sends message to Actor

Stream Processing – STORM

Core Concepts and Architecture

1. **Stream:**
 - o The core abstraction in Storm is the "stream," which is an unbounded sequence of tuples (data records).
 - o Streams are the primary input/output data format for all operations in a Storm topology.
2. **Tuple:**
 - o A tuple is a named list of values, essentially the smallest unit of data processed in a Storm application.
 - o Tuples can contain any number of fields, such as integers, strings, or objects.
3. **Topology:**
 - o A topology is a directed acyclic graph (DAG) that represents the flow of data between computations in Storm.
 - o A topology runs indefinitely until terminated, processing streams of data in real-time.
4. **Spouts:**
 - o Spouts are the sources of streams in a Storm topology.
 - o They emit data tuples into the topology, typically pulling data from external sources like message queues (e.g., Kafka), databases, or APIs.
5. **Bolts:**
 - o Bolts process tuples from spouts or other bolts. They are the building blocks for processing logic.
 - o Bolts can perform various tasks such as filtering, aggregation, joining streams, or interacting with databases.
6. **Stream Grouping:**
 - o Determines how tuples are routed between bolts.
 - o Examples:
 - **Shuffle Grouping:** Distributes tuples randomly.
 - **Fields Grouping:** Routes tuples to bolts based on specific fields.
 - **All Grouping:** Sends each tuple to all bolts.
 - **Global Grouping:** Sends all tuples to a single bolt.

Key Features

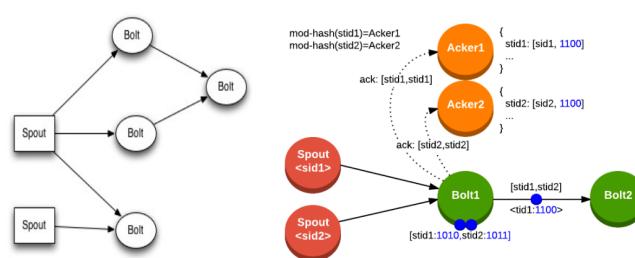
1. **Fault Tolerance:**
 - o Storm ensures reliability through an acknowledgment system. If a tuple fails to process, it can be replayed.
 - o Worker nodes can restart automatically in case of failure.
2. **Scalability:**
 - o Storm can handle a high throughput of messages by parallelizing computation.
 - o Parallelism is achieved by specifying the number of executors and tasks for spouts and bolts.
3. **Low Latency:**
 - o Designed for sub-second processing times, Storm is ideal for applications requiring immediate insights.
4. **Language Agnostic:**
 - o Storm's components can be written in any programming language that supports JSON serialization.

Stream Grouping

- Shuffle Grouping: Randomly
- Fields Grouping: Choose according to tuple field's value
- All Grouping: all tasks
- Global Grouping: send to tasks with the minimum id

Fault tolerance in Stream Processing

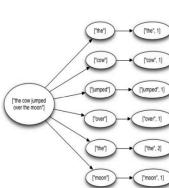
- Per element processing(at least once) - Storm
- Mini-batch(Exactly once) – SparkStream
- Per element processing(Exactly once) - Flink



1 Define the spout

Complete processing of a message

- When spouts emit tuples, specify a unique message ID.
- When bolts are done processing the input tuple, ack or fail the input tuple.



TestWordSpout

```
public void nextTuple() {
    Utils.sleep(100);
    final String[] words = new String[] {"nathan", "mike", "jackson", "golda", "bertels"};
    final Random rand = new Random();
    final String word = words[rand.nextInt(words.length)];
    _collector.emit(new Values(word));
}
```

emit() is to push the data to downstream bolts

WordCount in Storm

```
public static class WordCount implements IBasicBolt {
    Map<String, Integer> counts = new HashMap<String, Integer>();

    public void prepare(Map conf, TopologyContext context) {
    }

    public void execute(Tuple tuple, BasicOutputCollector collector) {
        String word = tuple.getString(0);
        Integer count = counts.get(word);
        if(count==null) count = 0;
        count++;
        counts.put(word, count);
        collector.emit(new Values(word, count));
    }

    public void cleanup() {
    }
}
```

Running Storm

```
LocalCluster cluster = new LocalCluster();

Map conf = new HashMap();
conf.put(Config.TOPOLOGY_DEBUG, true);

cluster.submitTopology("demo", conf, builder.createTopology());

cluster.shutdown();
```

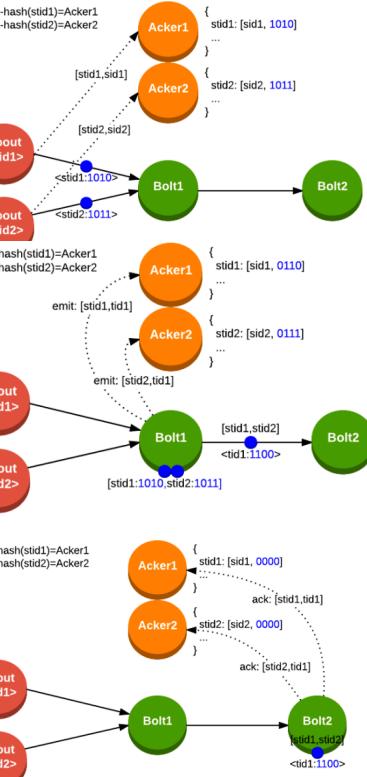
How to track if the message has been processed?

- We can follow the message flow to see if there are outstanding messages
- Use a data structure to record the generations/acknowledgement of messages
 - Do we need a new entry for each generated message?



Maintain the pending message

- Message 1 , acknowledged
 - Message 2, pending
 - ...
- erminate message 1 → Generate message 2
Process and acknowledge message 1



3. Define the topology

ExclamationTopology

```
TopologyBuilder builder = new TopologyBuilder();
builder.setSpout("words", new TestWordSpout(), 10);
builder.setBolt("exclaim1", new ExclamationBolt(), 1).shuffleGrouping("words");
builder.setBolt("exclaim2", new ExclamationBolt(), 2).shuffleGrouping("exclaim1");
```

2. Define the bolts

ExclamationBolt

```
public void execute(Tuple tuple) {
    _collector.emit(tuple, new Values(tuple.getString(0) + "!!!"));
    _collector.ack(tuple);
}

public void declareOutputFields(OutputFieldsDeclarer declarer) {
    declarer.declare(new Fields("word"));
}
```

Notice the ack() here, why is necessary?

SparkStream

Review: How Spark deal with Fault Tolerance

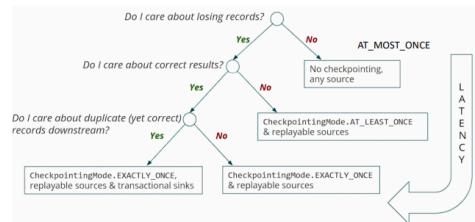
- Computation are implemented as deterministic transformations on immutable data sets

```
messages = textFile(...).filter(_.contains("error"))
           .map(_.split("\t")(2))
```



Different fault tolerance levels

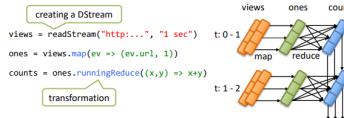
At most once, At least once, Exactly once



Example: Counting page views

Discretized Stream (DStream) is a sequence of immutable, partitioned datasets

- Can be created from live data streams or by applying bulk, parallel **transformations** on other DStreams



SparkStream WordCount code example*

```
// Create a ReceiverInputDStream on target ip:port and count the
// words in input stream of \n delimited test (e.g. generated by 'nc')
val lines = ssc.socketTextStream(args(0), args(1).toInt)
val words = lines.flatMap(_.split(" "))
val wordDStream = words.map(x => (x, 1))

// Update the cumulative count using mapWithState
// This will give a DStream made of state (which is the cumulative count of the words)
val mappingFunc = (word: String, one: Option[Int], state: State[Int]) => {
  val sum = one.getOrElse(0) + state.getOption.getOrElse(0)
  val output = (word, sum)
  state.update(sum)
  output
}
```

Discretized Stream Processing

Run a streaming computation as a series of small, deterministic batch jobs

Store intermediate state data in cluster memory

Try to make batch sizes as small as possible to get second-scale latencies

Fault tolerance in Stream Processing

- Per element processing(at least once) - Storm
- Mini-batch(Exactly once) – SparkStream
- Per element processing(Exactly once) - Flink

Discretized Stream Processing

Comparison of different Streaming Models

Comparison of different systems

- Programming model
 - Storm and Flink requires save and processing more
 - SparkStream is an extension of Spark
 - But Flink provide SQL like interface which is very user friendly

Flow control

- Push or pull
 - Spark uses push
 - Storm and Flink uses pull, buffer overflow only at the first stage
 - Pull is also known as back pressure

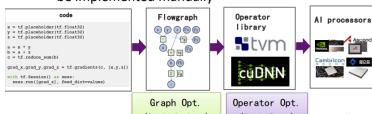
Fault tolerance

- Fault Tolerance
 - Storm guarantees a message is processed (at least) once
 - SparkStream uses mini-batch to provide exact-once
 - Flink (and storm trident) can do per-message exactly once with checkpoint, with replayable source(kafka) and transactional sink

Machine Learning

Machine Learning Systems

- Compile the program to flowgraph and optimize it
- Support libraries for new hardware
 - Thousands of operators which are not affordable to be implemented manually



Matrix Multiplication Loop Orders

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = 1.25

```
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        for (k=0; k<n; k++) {
            sum += a[i][j] * b[j][k];
        }
    }
}
```

kij (& ikj):

- 2 loads, 1 store
- misses/iter = 0.5

```
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        r = 0.0;
        for (k=0; k<n; k++) {
            r += a[i][k] * b[k][j];
        }
        c[i][j] = r;
    }
}
```

jki (& kji):

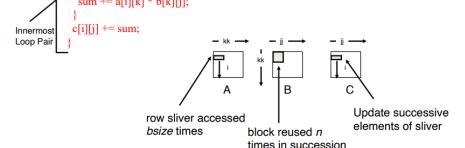
- 2 loads, 1 store
- misses/iter = 2.0

```
for (j=0; j<n; j++) {
    for (i=0; i<n; i++) {
        r = 0.0;
        for (k=0; k<n; k++) {
            r += a[j][k] * b[k][i];
        }
        c[j][i] = r;
    }
}
```

Blocked Matrix Multiply Analysis

- Innermost loop pair multiplies a $1 \times bsize$ sliver of A by a $bsize \times bsize$ block of B and accumulates into $1 \times bsize$ sliver of C
- Loop over i steps through n row slivers of A & C, using same B

```
for (i=0; i<n; i++) {
    for (j=jj; j< min(ij+bsize,n); j++) {
        sum = 0.0
        for (k=kk; k< min(kk+bsize,n); k++) {
            sum += a[i][k] * b[k][j];
        }
        c[i][j] += sum;
    }
}
```



Halide

Halide

- Halide is a domain specific language

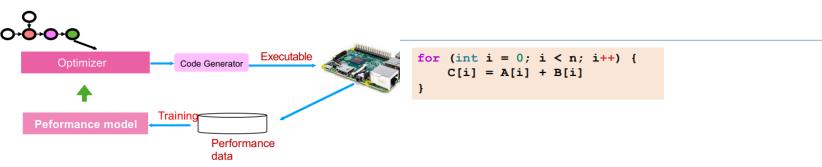
- It decoupled the optimized implementation

Program = algorithm + scheduling

- Enables programmers with different skills to collaborate
 - Algorithm experts and system experts
- Compilers are used to generate the parallel and optimized code

- Originally designed for image processing,
 - Now used by AI compilers for AI operator generation

Machine learning based optimizer



The machine learning model predict the performance within 1ms

Tensor-based abstraction

`C = matrix_multiply(A, B.T)`

```

m, n, h = tvm.var('m'), tvm.var('n'), tvm.var('h')
A = tvm.placeholder((m, h), name='A')
B = tvm.placeholder((n, h), name='B')

k = tvm.reduce_axis(0, h, name='k')
C = tvm.compute((m, n),
    lambda i, j: tvm.sum(A[i, k] * B[j, k], axis=k))
    
```

Inp
Output matrix shape
Computation rule:
$$C[i,j] = \sum_k A[i,k] \times B[j,k]$$

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

More tensor expressions

- Affine transformations

```

out = tvm.compute((n,m), lambda i, j: tvm.sum(data[i,k] * w[j,k], k))
out = tvm.compute((n,m), lambda i, j: out[i,j] + bias[i])
    
```

- Convolution

```

out = tvm.compute((c, h, w),
    lambda i, x, y: tvm.sum(data[kc,x+kx,y+ky] * w[i,kx,ky], [kx,ky]))
    
```

- ReLU

```

out = tvm.compute(shape, lambda *i: tvm.max(0, out(*i)))
    
```

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

```

xo, xi = s[C].split(s[C].axis[0], factor=32)
    
```

```

s[C].reorder(xi, xo)
    
```

```

s[C].bind(xo, tvm.thread_axis("blockIdx.x"))
    
```

```

s[C].bind(xi, tvm.thread_axis("threadIdx.x"))
    
```

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

```

xo, xi = s[C].split(s[C].axis[0], factor=32)
    
```

```

s[C].reorder(xi, xo)
    
```

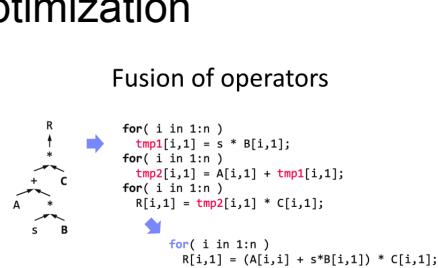
```

s[C].bind(xo, tvm.thread_axis("blockIdx.x"))
    
```

```

s[C].bind(xi, tvm.thread_axis("threadIdx.x"))
    
```

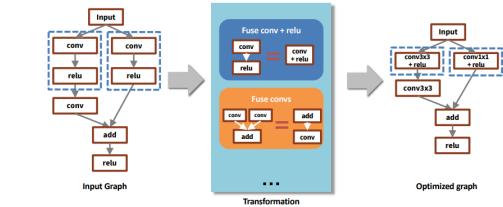
Transformation of Graphs



Fusion of operators

Graph Optimization

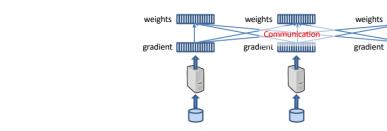
- Reduce the intermediate states which may cause memory references



Data Parallel Training

- In each iteration:

- each worker reads different mini-batch of size M, produces gradient
- communicate to produce aggregated gradient
- update weights

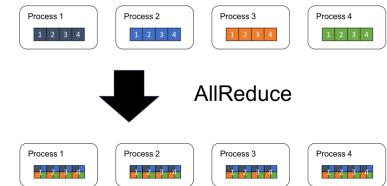
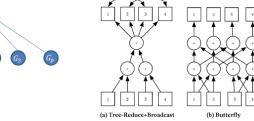


AllReduce to get global gradients

- We need to sum up all the gradients from different machines

$$G = G_1 + G_2 + \dots + G_p$$

- Ways to do the summation:

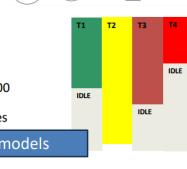
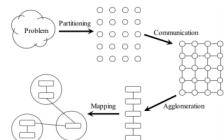


Distributed Training

Distributed training on 1000s cards

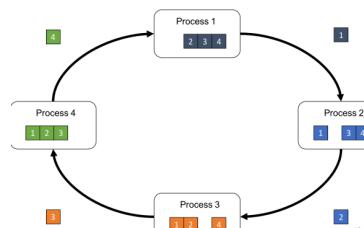
- Challenges
 - Communication reduction
 - Minimize the communications between between GPUs to allow high compute efficiency
 - Limited GPU memory
 - Intermediate state for backward computation
 - Load balance
 - Each card should have similar amount of work
 - Fault tolerance
 - MTBF is likely 10-100 hours, while training usually takes ~1000 hours
 - Quick recovery from system failures

10K – 100K cards for GPT4-like models

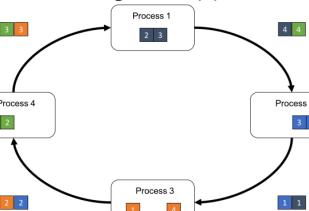


AllReduce to get global gradients

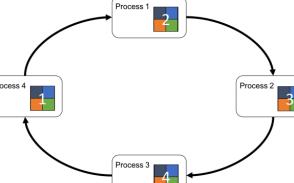
Ring Reduce(1)



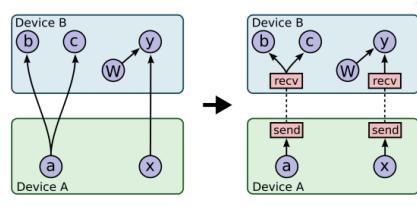
Ring Reduce(2)



Ring Reduce(3)

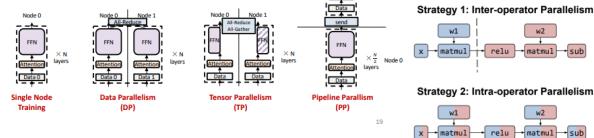


Model Parallel Training



Parallelisms for transformer

- Tensor model parallelism (split a layer)
- Pipeline parallelism (split between layers)

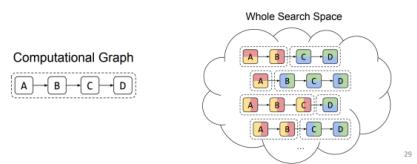


Challenges to computer systems :
How to partition the model and
choose parallelism?



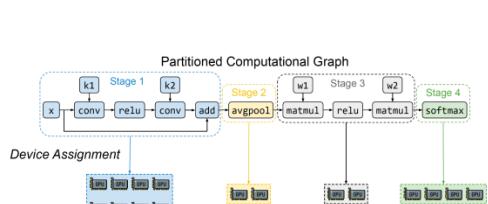
Can we find the good partition for a large model automatically

- Support both intra and inter-op parallelism, as well as pipeline parallelism
 - There are huge search space and it is a combinatorial optimization problem



Inter-op parallelism

Graph partition and device assignment within one compiler stage



Alpa

- Automatic parallelize the model training with just one line annotation

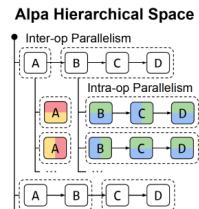
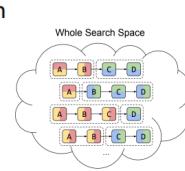
```

@alpa.parallelize
def train_step(model_state, batch):
    def loss_func(params):
        out = model_state.forward(params, batch["x"])
        return np.mean((out - batch["y"]) ** 2)

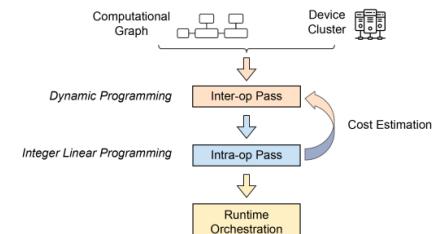
    grads = grad(loss_func)(state.params)
    new_model_state = model_state.apply_gradient(grads)
    return new_model_state

# A typical JAX training loop
model_state = create_train_state()
for batch in data_loader:
    model_state = train_step(model_state, batch)
    
```

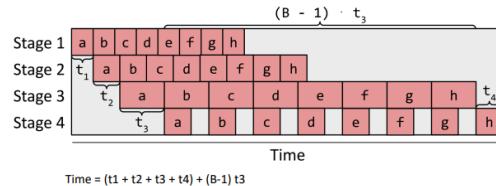
Two-level Hierarchical space to parallelism



Alpa Compiler phase



Inter-op parallelism



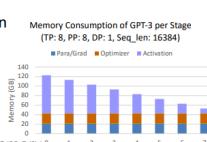
- Generalized objective function

$$T^* = \min_{\substack{s_1, \dots, s_S; \\ (n_1, m_1), \dots, (n_S, m_S)}} \left\{ \sum_{i=1}^S t_i + (B-1) \cdot \max_{1 \leq j \leq S} \{t_j\} \right\}$$

Memory Consumption

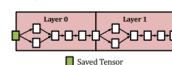
More on memory and load balance

- Large memory consumption
 - Optimizer
 - Parameters/Gradients
 - Activations
- Activations occupy huge memory
 - Recompute to save memory



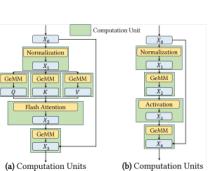
Recomputation

- Drop activations in the forward pass
- Recompute the activations in the backward pass
- Increased computation overhead**



Computation Units

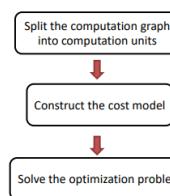
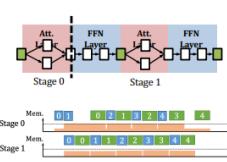
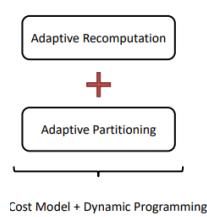
- Fine-grained partitioning
- Operators that require saving output tensors becomes the boundaries.
- Each Tensor bounds to one unit.



AdaPipe

AdaPipe

Adaptive Recomputation



Requires interfaces between each pair of entities

Blockchain

Motivation: Why correspondent banks are necessary— Lack of trust

Challenges of a global ledger

- Who can join this ledger?
 - Permissioned or permissionless
- Performance and availability
 - Global network latency and network partition
- Who can write what on this ledger?
 - How to avoid fake transactions?
 - How to avoid double spend?
 - How to avoid modification to previous transactions?
- Who should provide computers to maintain this ledger?
 - Why are they interested in doing this - Incentive

Who can join the ledger?

- Permissionless
 - Public chain
- Permissioned
 - Within a single organization
 - Private chain
 - Multiple party
 - Consortium Blockchain

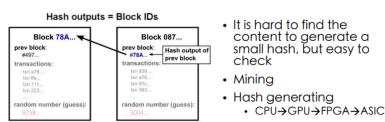


Immutable data



- Crypto-hashing
 - Sha256
 - Single way function, sha256(byte array) = 256bits hash
 - Hard to construct byte array with a given hash
 - A block id which is the hash of the transactions in the block
 - A block contains the blockid of the previous block

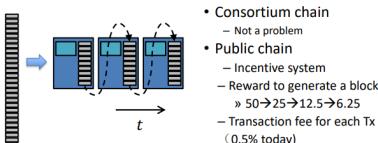
Proof of work



Choose transactions that:

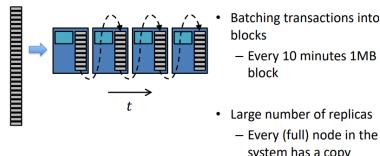
$\text{Sha256}(\text{prevID}, \text{transactions}, \text{guess}) < 00\ldots00010000$

(4) Why people are volunteering to maintain the blockchain systems

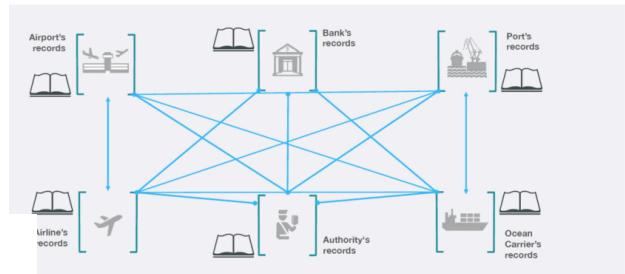


- Consortium chain
 - Not a problem
- Public chain
 - Incentive system
 - Reward to generate a block
 - » 50 → 25 → 12.5 → 6.25
 - Transaction fee for each Tx (0.5% today)

Performance and availability



- Batching transactions into blocks
 - Every 10 minutes 1MB block
- Large number of replicas
 - Every (full) node in the system has a copy



RSA (Rivest–Shamir–Adleman)

It is commonly utilized to ensure secure communication and for creating digital signatures. It Uses large integer prime numbers for key generation. It Encrypts data with the public key and decrypts with the private key. It is slower than some other algorithms but offers strong security.

Key Generation

- Choose>Select two large prime numbers, p and q.
- Calculate $n = p \cdot q$.
- Calculate $\phi(n) = (p-1)(q-1)$, where ϕ is Euler's totient function.
- Choose an integer e, that $1 < e < \phi(n)$ and $\gcd(e, \phi(n)) = 1$.
- Compute d, the modular multiplicative inverse of e modulo $\phi(n)$.
- Public key: (e, n)
- Private key: (d, n)

Encryption

- Convert plaintext message into an integer m.
- Compute ciphertext $c = m^e \bmod n$.

Decryption

- Calculate the plaintext message $m = c^d \bmod n$.

Select two prime no's. Suppose $P = 53$ and $Q = 59$.

Now First part of the Public key : $n = P \cdot Q = 3127$.

We also need a small exponent say e :

But e Must be an integer. Not be a factor of $\Phi(n)$.

$1 < e < \Phi(n)$ $[\Phi(n) = (P-1)(Q-1) = 3016]$,

Choose $e = 3$. public key $(e, n) = (3, 3127)$

$d = (k^*\Phi(n) + 1) / e$ for some integer k

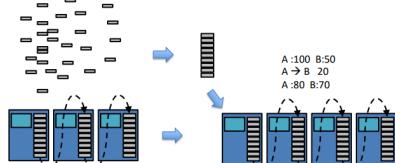
For $k = 2$, value of d is 2011.

Private key $(d, n) = (2011, 3127)$

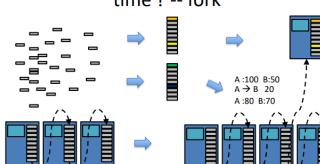
Encrypted Data $c = (89^e) \bmod n = 1394$

Decrypted Data $= (c^d) \bmod n = 89$

Consensus : Agree on the content of transactions in the block



What happens if two of the people solve the challenges at (almost) same time ? -- fork



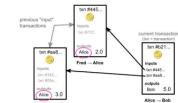
The Longest chain rule

Bitcoin blockchain

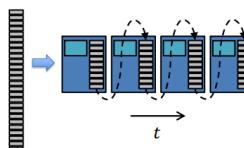
- Immutable, all transactions are pu

How to know if a transaction is Unspent?

- Track the whole block chain?
 - 873,585 today
 - A large data structure contains many spent output
- A UTXO database
 - Only unspent transactions



Put it together : The bitcoin blockchain



- A system which integrates cryptography, distributed systems, game theory, monetary finance
- Code is law
- Infrastructure for trust

Challenges with block chains

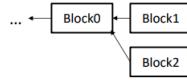
- Public chain
 - Poor scalability
 - Slow (Bitcoin: 7 TPS, Taobao, VISA: 100K TPS(peak), UBER: 100TPS)
- Consortium chain
 - Poor implementation
- Security
 - Both on systems and smart contract
- Privacy
 - Zero knowledge proof and Zcash
 - Content

Problems of block chain

Forks in Nakamoto Consensus

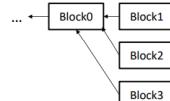
Scalability:

What if we run Nakamoto Consensus with **large/fast** blocks generation?



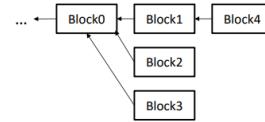
Large blocks → Longer propagation time
Fast block generation → Less time for propagation

Forks in Nakamoto Consensus



More forks when the block generation is fast and

Forks in Nakamoto Consensus



Usually the system eventually converges on one fork by waiting for more blocks appended in the end

Solution:



Projects:



Wait! Why do we blockchain in the first place?

Great idea. But need large players to deposit on the chain.

Centralize



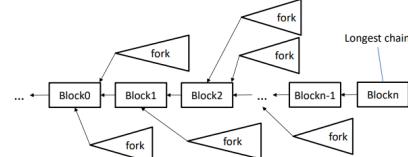
Off-chain

Against CAP Theorem -- rely on economical mechanism to secure each shard.

Sharding



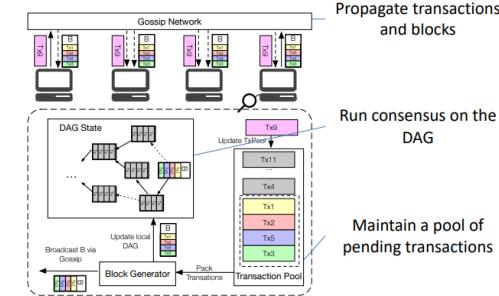
Forks in Nakamoto Consensus



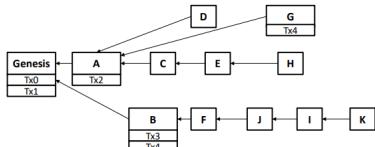
- Forks waste network/processing resources
- Downgrade safety
- Attacker needs less resource to beat the longest chain

Conflux

Conflux Architecture



Tx0: Mint 10 coins to X
Tx1: Mint 10 coins to Y
Tx2: X sends 8 to Y
Tx3: X sends 8 to Z
Tx4: Y sends 8 to Z

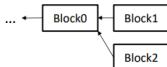


Each block has one outgoing parent edge to its parent block

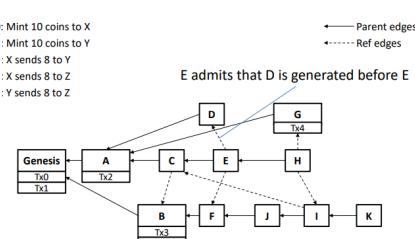
Parent edges form a tree

One Observation

- Bitcoin enforces a restrictive transaction total order at the generation time of each block:



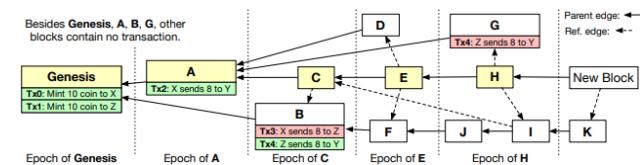
- Blockchain transactions rarely conflict with each other and they can be serialized in any order
- Why not process non-conflicting transactions in all concurrent blocks?



We still need a total order, how?

Conflux Consensus Ideas

- Optimistically process concurrent blocks without discarding any for forks
- Organize blocks into a **direct acyclic graphs (DAG)**
- Enable **fast/large block generation** without sacrificing security
 - First agree on a total order of all blocks
 - Assume transactions would not conflict each other
 - Then derives the transaction order from the agreed block order
 - Resolve transaction conflicts lazily



Computer Systems

Why is it difficult to build a system?

Scales of Computer Systems

- Programming / Data Structure
 - LoC (Lines of Code): From hundreds to millions
- Operating System / Architecture
 - Cores: from one to hundreds
- Network
 - Nodes: from two to millions
- Web Service
 - Clients: from tens to millions

Complexity of Computer Systems

- Hard to define; symptoms:
 - Large number of components
 - Large number of connections
 - Irregular
 - No short description
 - Many people required to design/maintain
- Technology rarely the limit
 - Indeed tech opportunity is the problem
 - Limit is usually designers' understanding

Common Problems of Systems

Problem Types

- **Emergent properties** (surprise!)
 - The properties that are not considered at design time
- **Propagation of effects**
 - Small change -> big effect
- **Incommensurate scaling**
 - Design for small model may not scale
- **Trade-offs**
 - Waterbed effect

Waterbed Effect

- Pushing down on a problem at one point
- Causes another problem to pop up somewhere else

Coping with Complexity

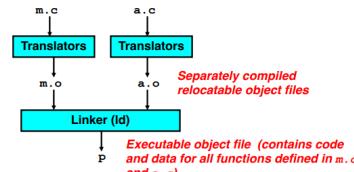
M.A.L.H

- **Modularity**
 - Split up system
 - Consider separately
- **Layering**
 - Gradually build up capabilities
- **Abstraction**
 - Interface/Hiding
 - Avoid propagation of effects
- **Hierarchy**
 - Reduce connections
 - Divide-and-conquer

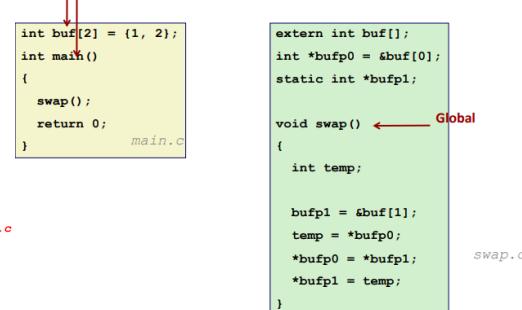
Layering

- Goal
 - Reduce module interconnections even more
- How to do it
 - Build a set of mechanisms first(a lower layer)
 - Use them to create a different complete set of mechanisms (an upper layer)
- General rule: A module in one layer only interacts with:
 - Its peers in the same layer, and
 - Modules in the next lower layer / next higher layer

Use linker to support modularity in software design



How pointers break modular boundaries



Unbounded Composition

- Two properties of computer systems
 - 1. Mostly digital
 - 2. Controlled by software
- Both relax the limits on complexity arising from physical laws in other systems

Prep. Exam

1. (15 points) In a storm application, a Spout with <sid1> sends a tuple to Bolt1 whose stid is 0010. It's acker is Ack1. After the execution of Bolt1, it sends two tuples to Bolt2 and Bolt3 respectively.

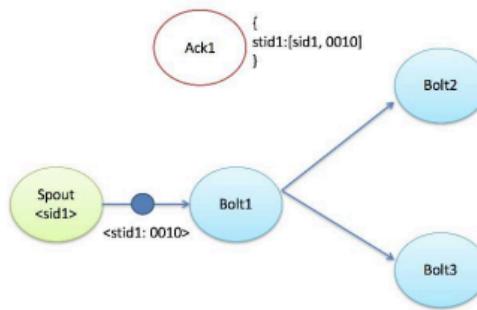
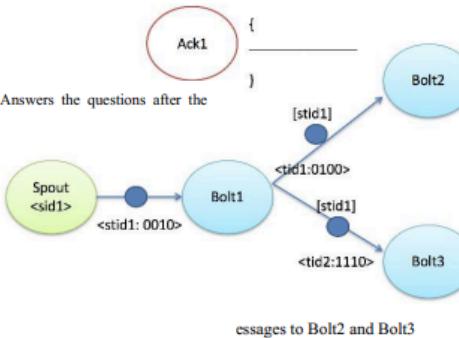


Figure 1.1 Spout sid1 sends a tuple to Bolt 1

3. (15points) For the following MPI algorithm. Answers the questions after the source code.



messages to Bolt2 and Bolt3

```

#include <stdlib.h>
#include <stdio.h>
#include <math.h>
#include <time.h>
#include "mpi.h"
int main ( int argc, char *argv[] );
int prime_number ( int n, int id, int p );

int main ( int argc, char *argv[] )
{
    int i, id, ierr, n, n_factor, n_hi, n_lo, p, primes, primes_part;
    double wtime;
    n_lo = 1;
    n_hi = 262144;
    n_factor = 2;
    ierr = MPI_Init ( &argc, &argv );
    ierr = MPI_Comm_size ( MPI_COMM_WORLD, &p );
    ierr = MPI_Comm_rank ( MPI_COMM_WORLD, &id );
    ierr = MPI_Wtime ( &wtime );
    primes_part = prime_number ( n, id, p );
    ierr = MPI_Reduce ( &primes_part, &primes, 1, MPI_INT,
    MPI_SUM, 0, MPI_COMM_WORLD );
    if ( id == 0 ) {
        wtime = MPI_Wtime ( ) - wtime;
        printf ( "%8d %8d %14fn", n, primes, wtime );
    }
    ierr = MPI_Finalize ( );
    return 0;
}

int prime_number ( int n, int id, int p )
{
    Purpose:
        PRIME_NUMBER returns the number of primes between 1 and N.
    Parameters:
        Input, int N, the maximum number to check.
        Input, int ID, the ID of this process,
        between 0 and P-1.
        Input, int P, the number of processes.
        Output, int PRIME_NUMBER, the number of prime numbers up to N.
    {
        int i, total, prime;
        total = 0;

        for ( i = 2 + id; i <= n; i = i + p )
            prime = 1;
            for ( j = 2; j < i; j++ )
                if ( ( i % j ) == 0 )
                    prime = 0;
                    break;
                }
            total = total + prime;
            }
        return total;
}
  
```

2. (15 points) In graph computing frameworks, sometimes it is useful to sort vertex with the degree, such as the bottom-up algorithm for BFS. Now we need to write a spark program for this task.

Input:

An unsorted edge list of a directed graph with format(source_vertex_id, destination_vertex_id), the vertex id is a unsigned integer less than 4 billion. The edge list is in text format, one edge per line. For example.

1,2
2,4
3,5
1,4

Output:

Vertex id sorted by in-degree with descending order. If two vertices have the same degree, they are sorted by vertex id with ascending order. The output of the above example is:

4
2
5
1
3

Question:

Write a spark program to perform the task.

```

lines = sc.textFile(input)
edges = lines.map(lambda line: tuple(line.split(',')))
degree_v = edges.map(lambda e: (e[1], 1)).reduce(add).sortBy(lambda x: (-x[1],
x[0])).collect()
for _d, v in degree_v:
    print(v)
  
```

- a) Please fill the blanks in the Figure 1.2 to show the record of stid1. The format example is shown in Figure 1.1.
After Bolt1 send ack to Ack1, record will be {stid1: [stid1: 1010]}
- b) How would the system detect that a message is not processed completely and need to be reprocessed?
If a bolt crashes, or emits a failure not ack, or the tuple reaches timeout.
- c) What is the feature of the Spout to enable the feature of message reprocessing? What is the typical software used to support this feature?
Replayable source. Queue (Kafka).

Question 1: Fill in the blanks in the first page.

Question 2: Is this a load balanced algorithm? How many numbers are processed by each process with number n and process number p(roughly) in one iteration?

Yes. (n-2)/p.

Question3: Any ideas to optimize this code?

Many ways. For example, for (j = 2; j < i; j++) -> for (j = 2; j*j < i; j++).

